



# CS 4803 / 7643: Deep Learning

## Topics:

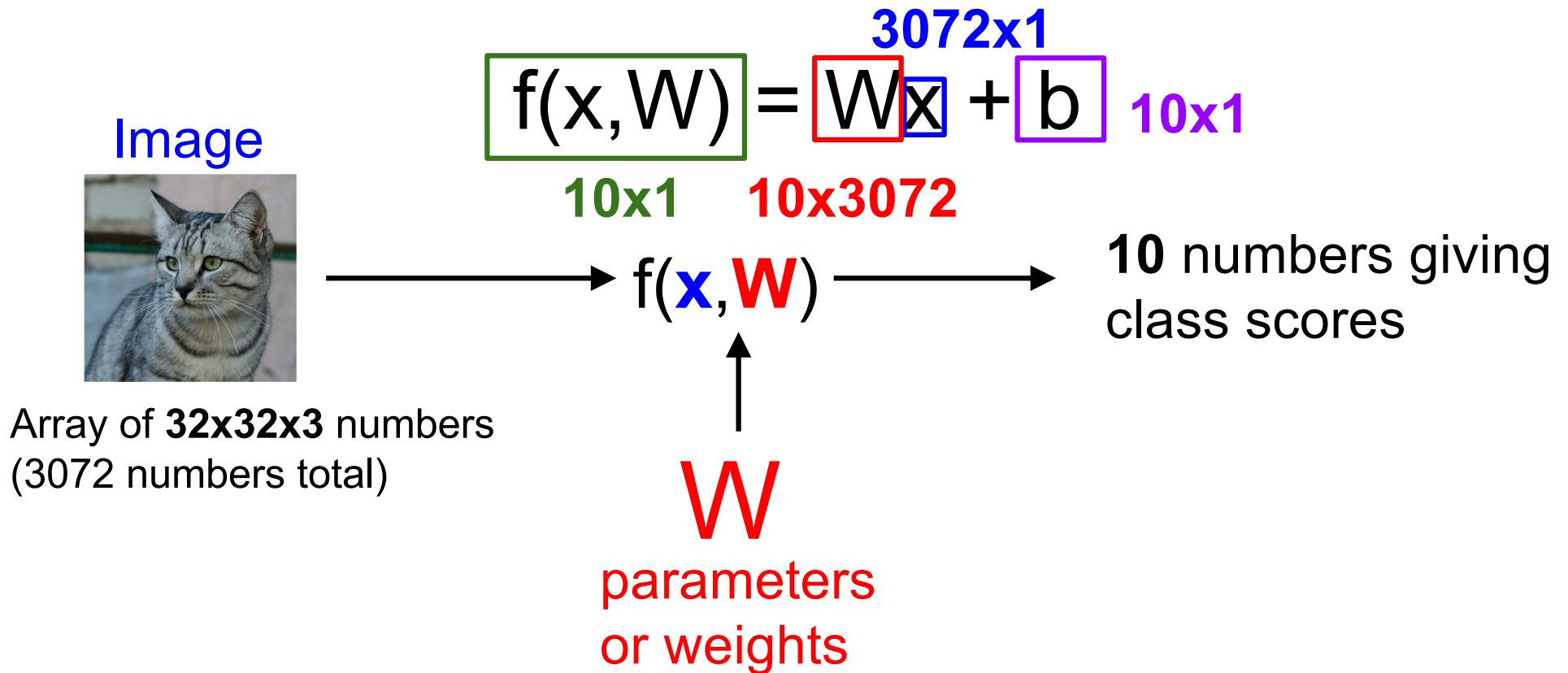
- Regularization
- Neural Networks
- Optimization
- Computing Gradients

Dhruv Batra  
Georgia Tech

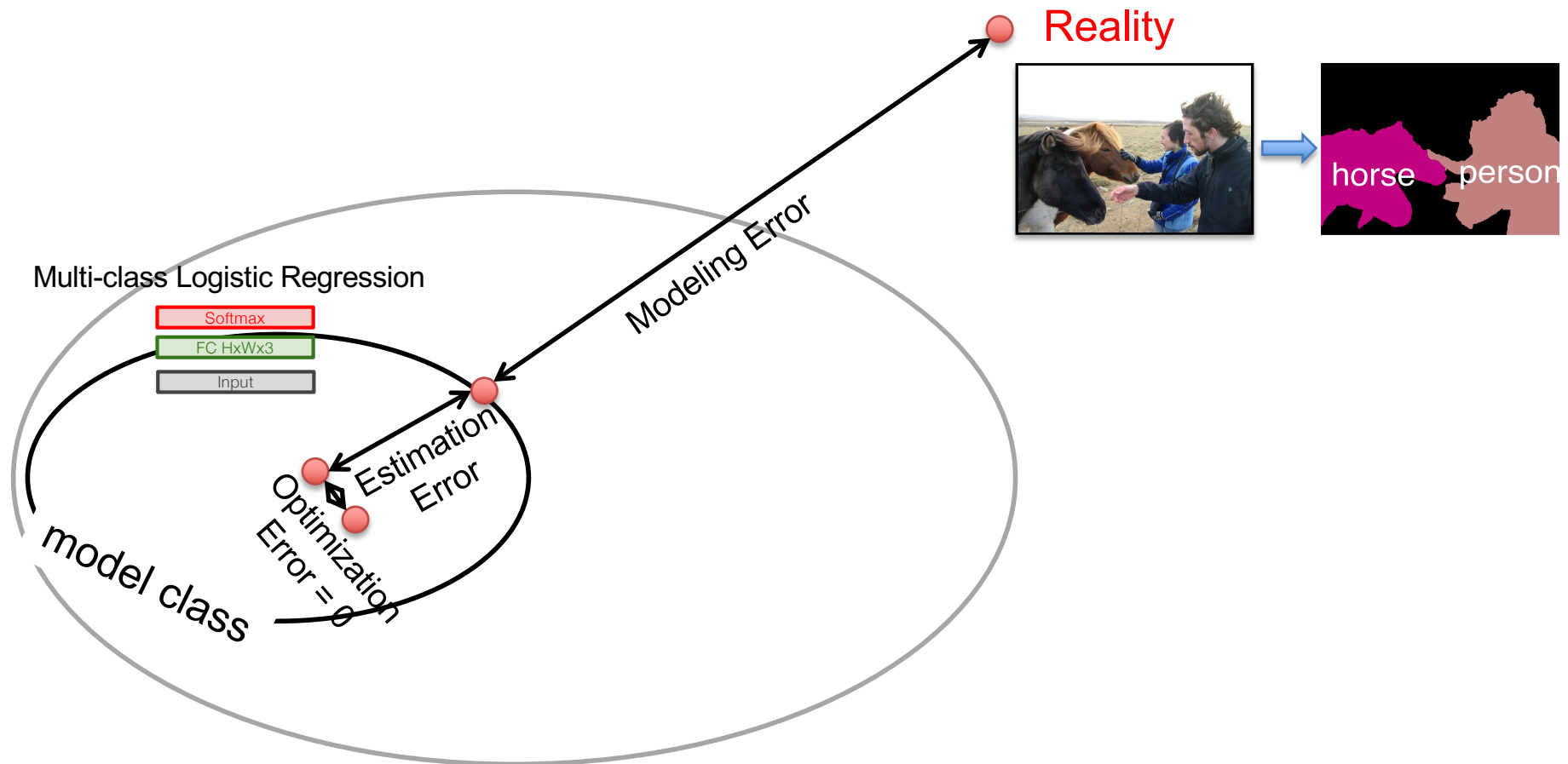


# Recap from last time

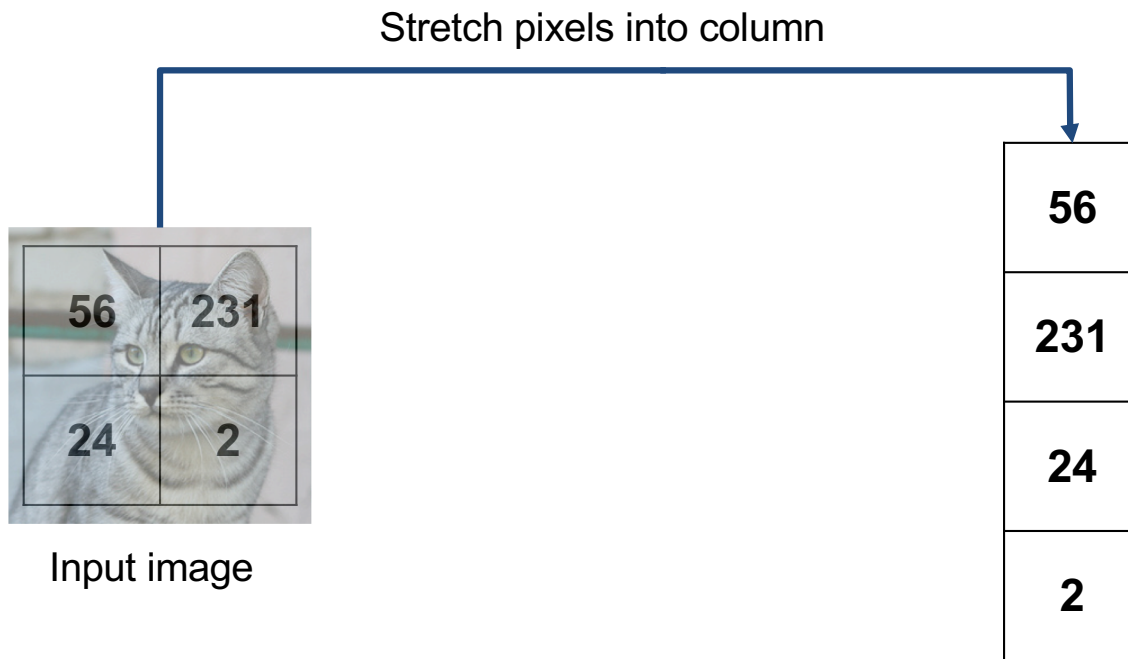
# Parametric Approach: Linear Classifier



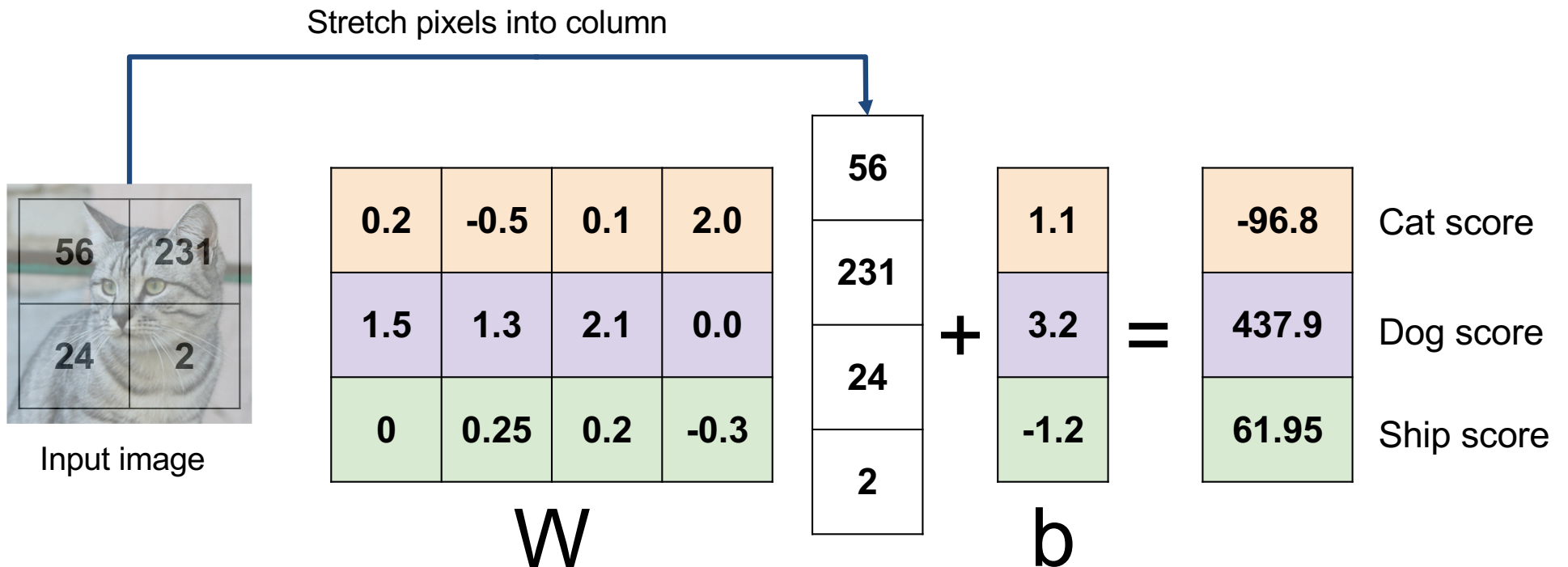
# Error Decomposition



Example with an image with 4 pixels, and 3 classes (cat/dog/ship)



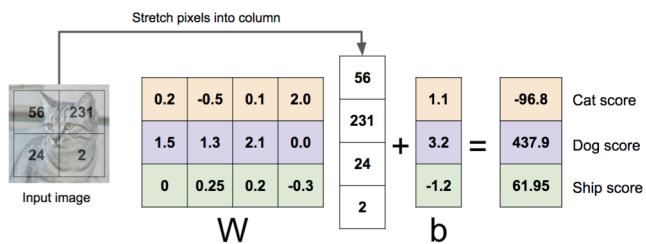
Example with an image with 4 pixels, and 3 classes (cat/dog/ship)



# Linear Classifier: Three Viewpoints

## Algebraic Viewpoint

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



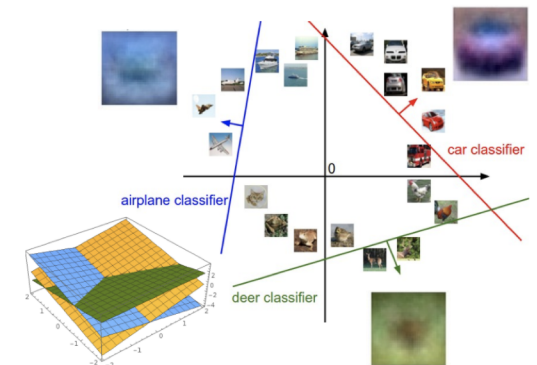
## Visual Viewpoint

One template  
per class



## Geometric Viewpoint

Hyperplanes  
cutting up space



# Recall from last time: Linear Classifier



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

[Cat image](#) by [Nikita](#) is licensed under [CC-BY 2.0](#); [Car image](#) is [CC0 1.0](#) public domain; [Frog image](#) is in the public domain

## 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

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

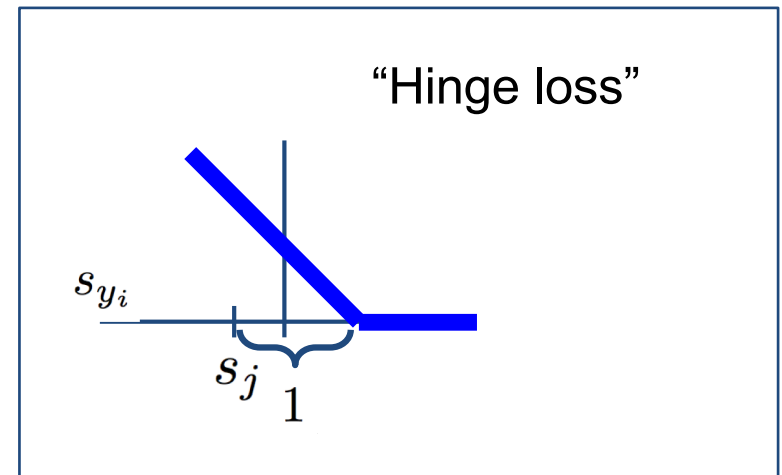
$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Suppose: 3 training examples, 3 classes.  
 With some  $W$  the scores  $f(x, W) = Wx$  are:



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>

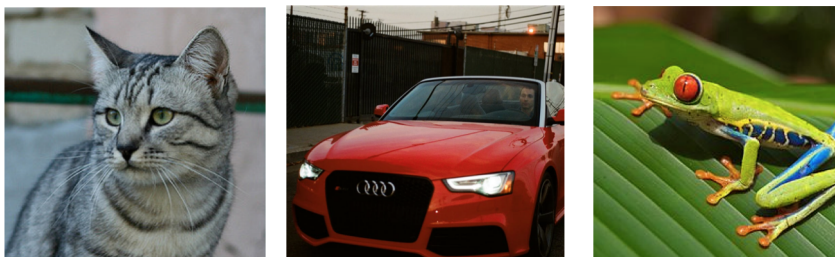
### Multiclass SVM loss:



$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

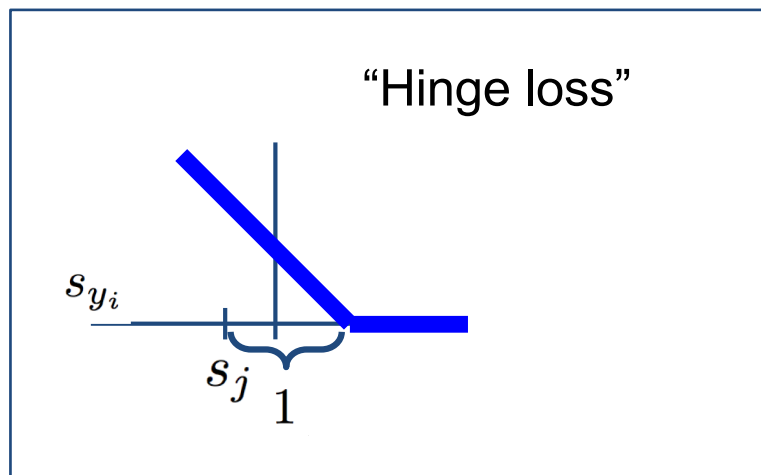
$$= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Suppose: 3 training examples, 3 classes.  
 With some  $W$  the scores  $f(x, W) = Wx$  are:



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>

### Multiclass SVM loss:



$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

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



## Softmax vs. SVM

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

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

## Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax  
Function

cat	<b>3.2</b>
car	5.1
frog	-1.7

# Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax  
Function

Probabilities  
must be  $\geq 0$

cat	3.2	exp →	24.5
car	5.1		164.0
frog	-1.7		0.18

unnormalized probabilities

# Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax  
Function

Probabilities  
must be  $\geq 0$

Probabilities  
must sum to 1

cat	3.2
car	5.1
frog	-1.7

exp

24.5
164.0
0.18

normalize

0.13
0.87
0.00

unnormalized  
probabilities

probabilities

# Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax  
Function

Probabilities  
must be  $\geq 0$

Probabilities  
must sum to 1

cat  
car  
frog

3.2
5.1
-1.7

Unnormalized log-  
probabilities / logits

exp

24.5
164.0
0.18

unnormalized  
probabilities

normalize

0.13
0.87
0.00

probabilities



# Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad \text{Softmax Function}$$

Probabilities  
must be  $\geq 0$

Probabilities  
must sum to 1

$$L_i = -\log P(Y = y_i|X = x_i)$$

cat  
car  
frog

3.2  
5.1  
-1.7

Unnormalized log-  
probabilities / logits

exp

24.5  
164.0  
0.18

unnormalized  
probabilities

normalize

0.13  
0.87  
0.00

probabilities

$$\rightarrow L_i = -\log(0.13) = 2.04$$

# Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad \text{Softmax Function}$$

Probabilities  
must be  $\geq 0$

Probabilities  
must sum to 1

$$L_i = -\log P(Y = y_i | X = x_i)$$

cat  
car  
frog

3.2  
5.1  
-1.7

Unnormalized log-  
probabilities / logits

exp

24.5  
164.0  
0.18

unnormalized  
probabilities

normalize

0.13  
0.87  
0.00

probabilities

$$\rightarrow L_i = -\log(0.13) = 2.04$$

**Maximum Likelihood Estimation**  
Choose probabilities to maximize  
the likelihood of the observed data

# Log-Likelihood / KL-Divergence / Cross-Entropy

$$D = \{(x_i, y_i)\}$$

$$\text{IID} \sim P^x$$

$$\begin{aligned} \hat{w}_{MLE} &= \max_w P(D|w) \\ &\equiv \max_w \log P(D|w) \\ &\equiv \max_w \sum_i \log P(y_i|x_i, w) \end{aligned}$$

# Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad \text{Softmax Function}$$

Probabilities must be  $\geq 0$

Probabilities must sum to 1

$$L_i = -\log P(Y = y_i | X = x_i)$$

cat  
car  
frog

3.2
5.1
-1.7

Unnormalized log-probabilities / logits

exp

24.5
164.0
0.18

unnormalized probabilities

normalize

0.13
0.87
0.00

probabilities

compare

1.00
0.00
0.00

Correct probs

# Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax Function

Probabilities must be  $\geq 0$

Probabilities must sum to 1

$$L_i = -\log P(Y = y_i | X = x_i)$$

cat  
car  
frog

3.2
5.1
-1.7

Unnormalized log-probabilities / logits

exp

24.5
164.0
0.18

unnormalized probabilities

normalize

0.13
0.87
0.00

probabilities

compare

Kullback-Leibler divergence

$$D_{KL}(P||Q) = \sum_y P(y) \log \frac{P(y)}{Q(y)}$$

1.00
0.00
0.00

Correct probs

# Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad \text{Softmax Function}$$

Probabilities must be  $\geq 0$

Probabilities must sum to 1

$$L_i = -\log P(Y = y_i | X = x_i)$$

cat  
car  
frog

3.2
5.1
-1.7

Unnormalized log-probabilities / logits

exp

24.5
164.0
0.18

unnormalized probabilities

normalize

0.13
0.87
0.00

probabilities

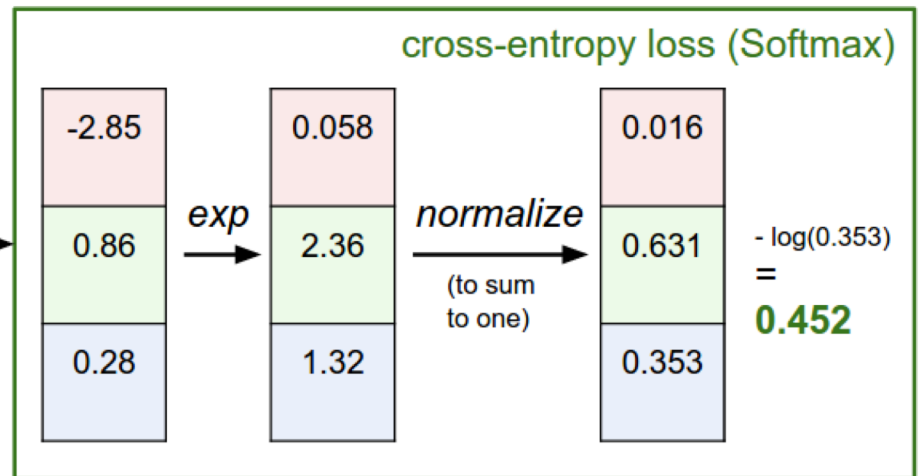
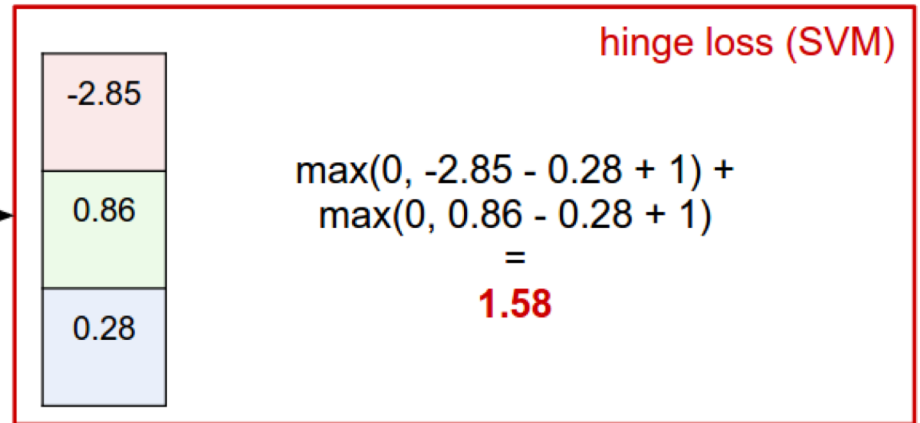
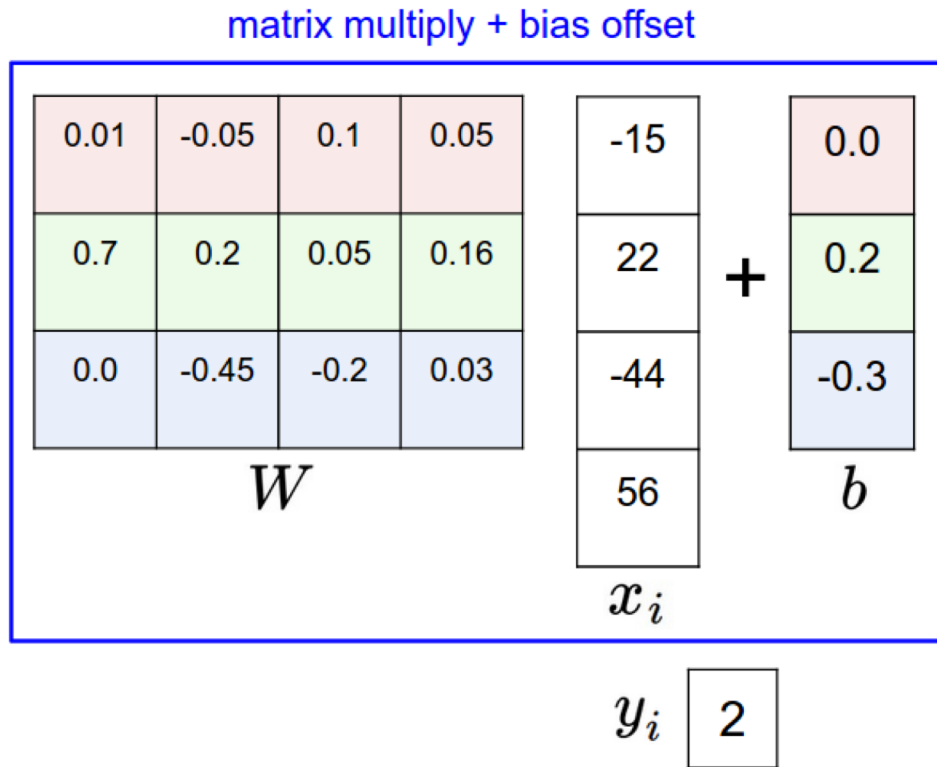
compare

1.00
0.00
0.00

Correct probs

Cross Entropy

$$H(P, Q) = H(p) + D_{KL}(P || Q)$$




# Plan for Today

- Regularization
- Neural Networks
- Optimization
- Computing Gradients




# Regularization

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


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

# Regularization


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


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

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

# Regularization

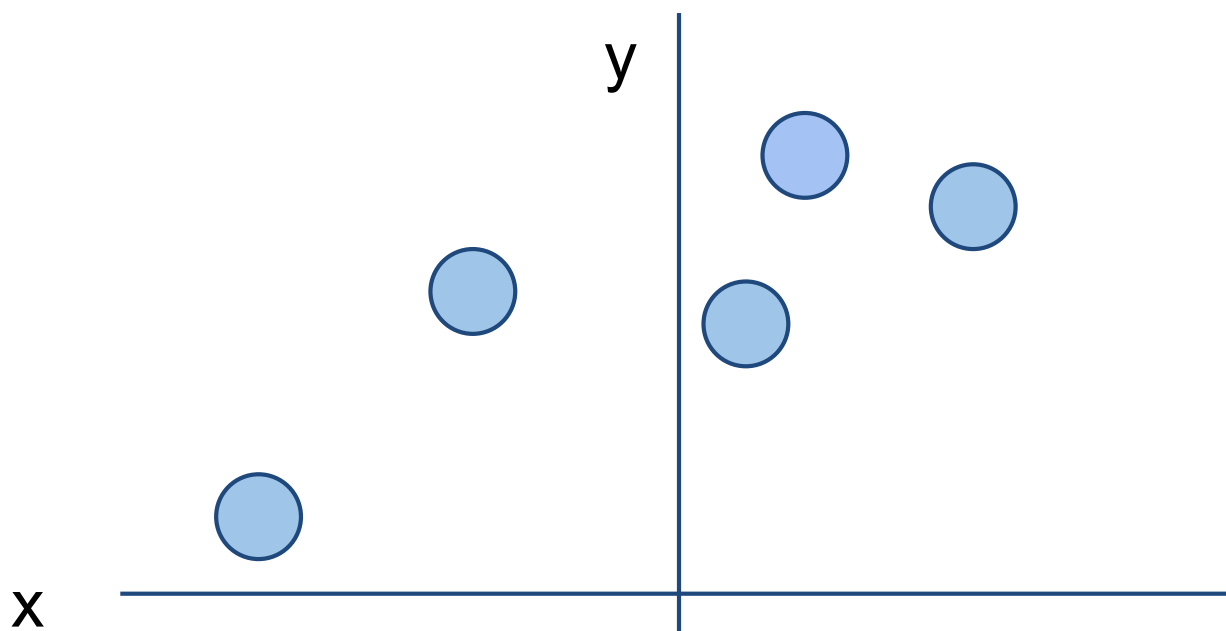
$\lambda$  = regularization strength  
(hyperparameter)

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


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

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

# Regularization: Prefer Simpler Models



# Polynomial Regression

$$\hat{y} = w_0 + w_1 x$$

$$\hat{y} = w_0 + w_1 x + w_2 x^2 + \dots + w_d x^d$$

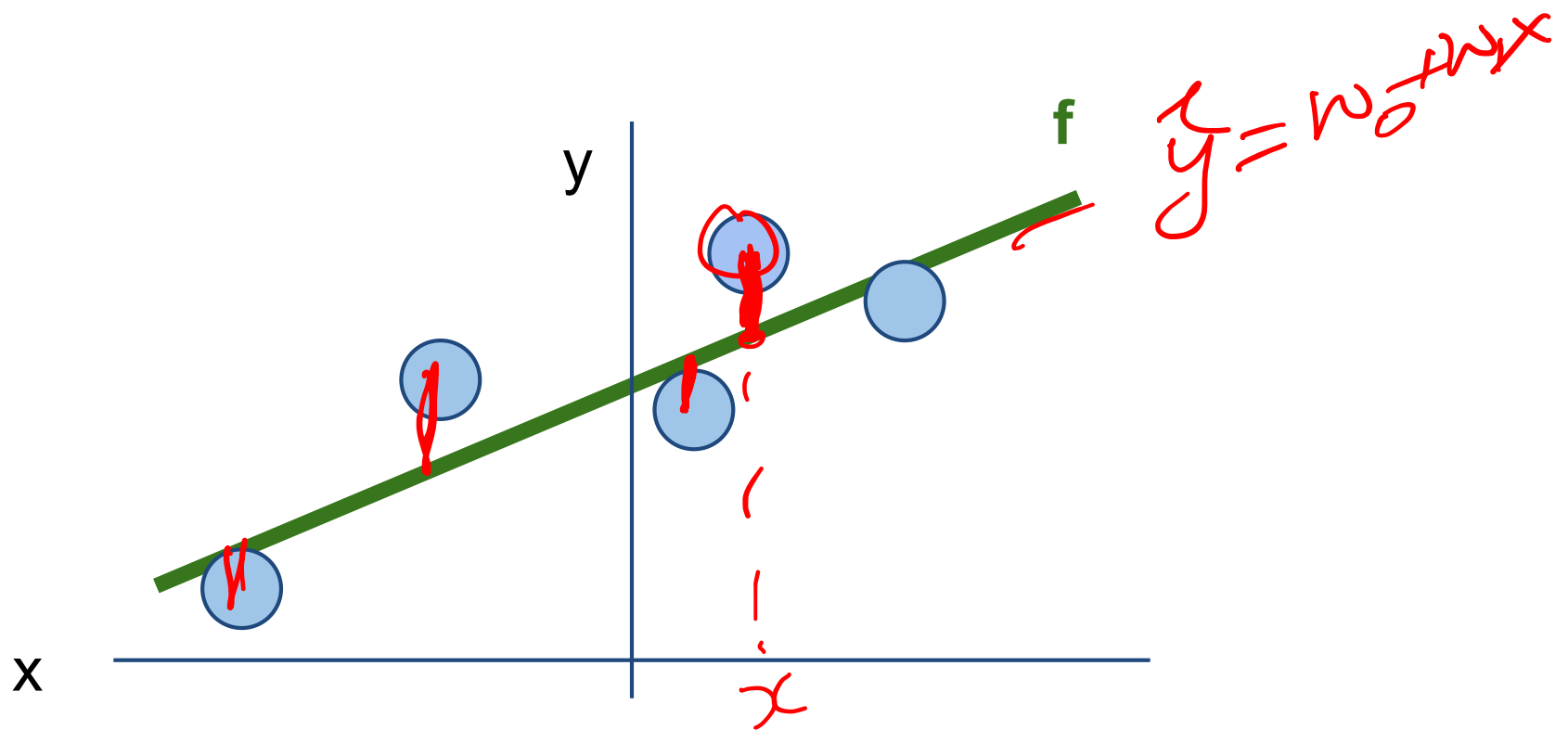
$$= [w_0 \dots w_d] \underbrace{\begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^d \end{bmatrix}}_{\phi(x)} = \vec{w}^T \phi(x)$$

$$D = \{(x_i, y_i)\}$$

$$w^* = \min_w \frac{1}{2} \sum (y_i - \hat{y}_i)^2$$

# Polynomial Regression

# Polynomial Regression




# Polynomial Regression

- Demo:
  - <https://arachnoid.com/polysolve/>
- Data:
  - 10 6
  - 15 9
  - 20 11
  - 25 12
  - 29 13
  - 40 11
  - 50 10
  - 60 9



# 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

# Regularization

$$L(W) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)}_{\text{Data loss: Model predictions should match training data}} + \underbrace{\lambda R(W)}_{\text{Regularization: Prevent the model from doing too well on training data}}$$

$\lambda$  = regularization strength (hyperparameter)

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

# 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

# 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

Why regularize?

- Express preferences over weights
- Make the model *simple* so it works on test data
- Improve optimization by adding curvature

$$\hat{y} = w_0 + w_1 x + \dots + \boxed{w_d} x^d$$

# Recap

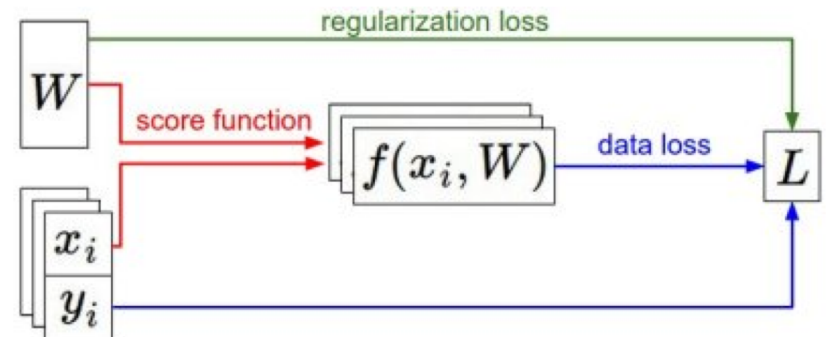
- We have some dataset of  $(x, y)$
- We have a **score function**:  $s = f(x; W) \stackrel{\text{e.g.}}{=} Wx$
- We have a **loss function**:

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

SVM

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

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



# Recap

How do we find the best  $W$ ?

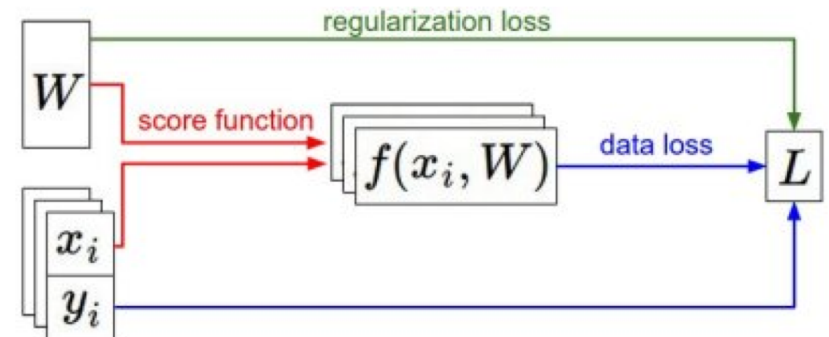
- We have some dataset of  $(x, y)$
- We have a **score function**:  $s = f(x; W) \stackrel{\text{e.g.}}{=} Wx$
- We have a **loss function**:

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

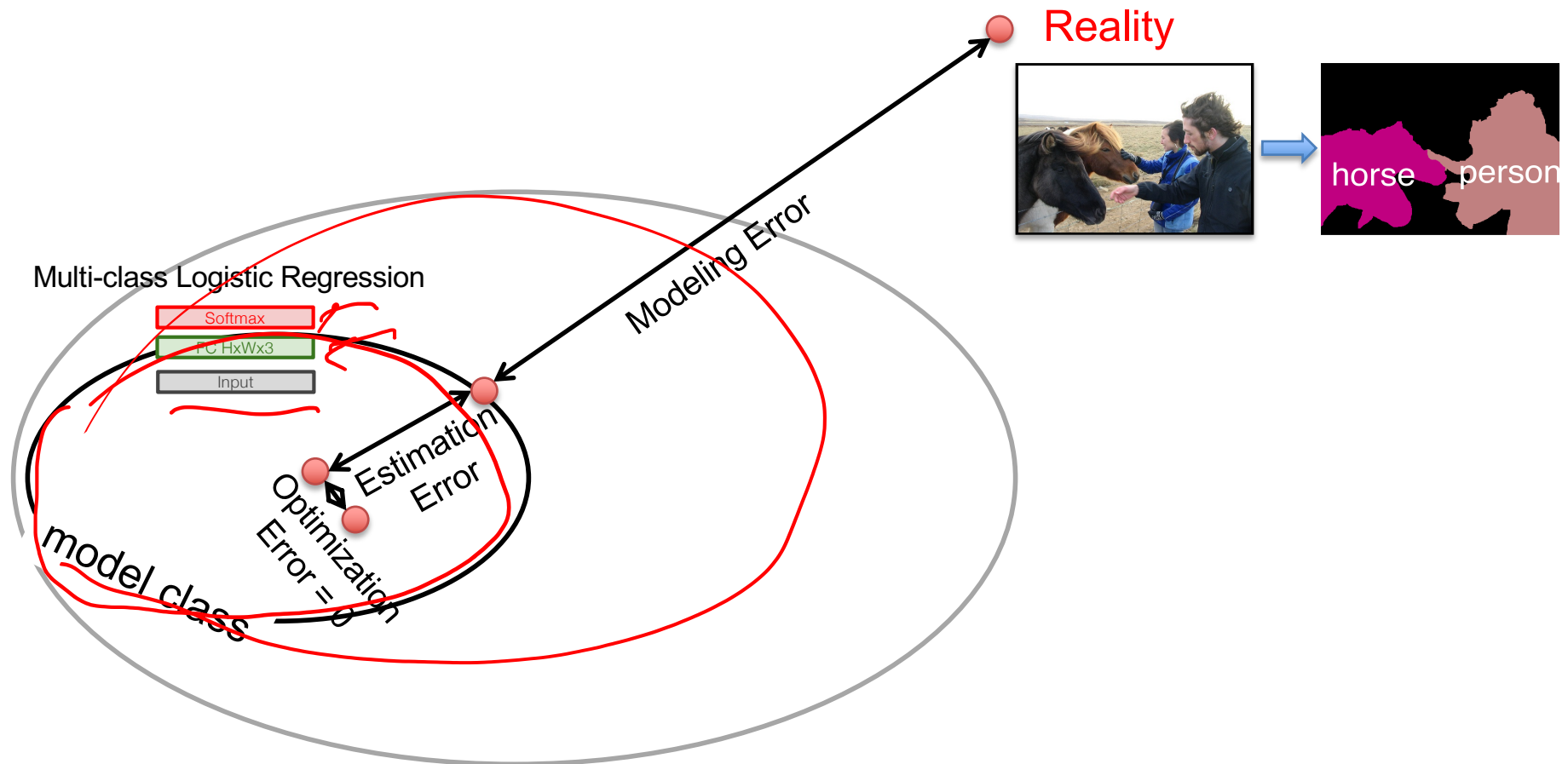
SVM

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

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



# Error Decomposition





# Next: Neural Networks



# Neural networks: without the brain stuff

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

$$f = \underline{w_2} w_1 x$$

# Neural networks: without the brain stuff

(**Before**) Linear score function:

$$f = Wx$$

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

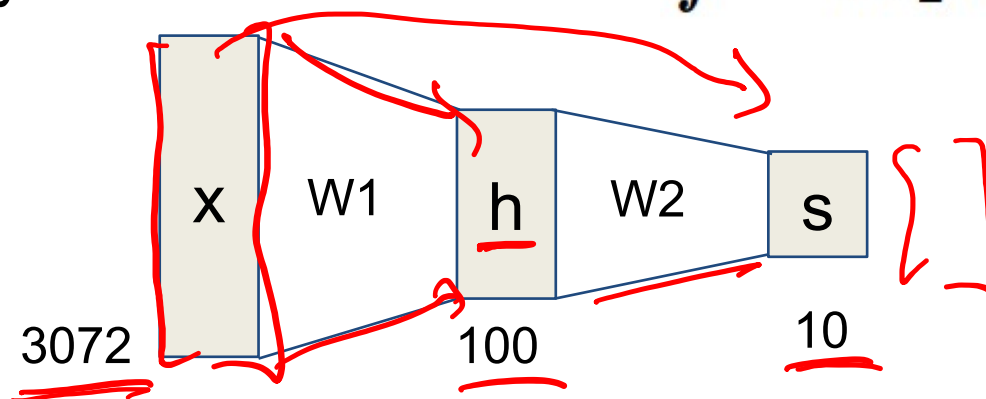
$$f = W_2 \max(0, W_1 x)$$



# Neural networks: without the brain stuff

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

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



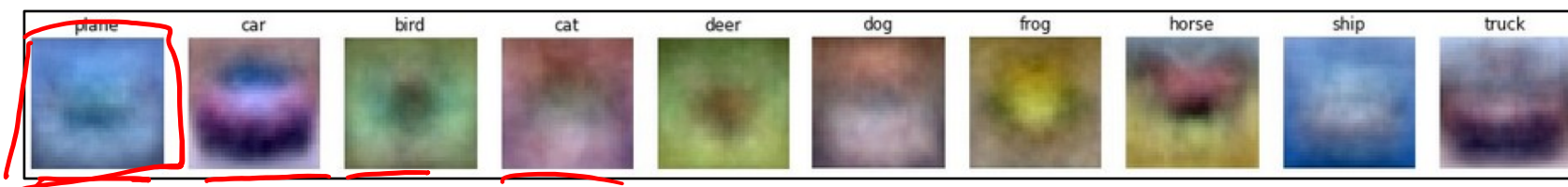
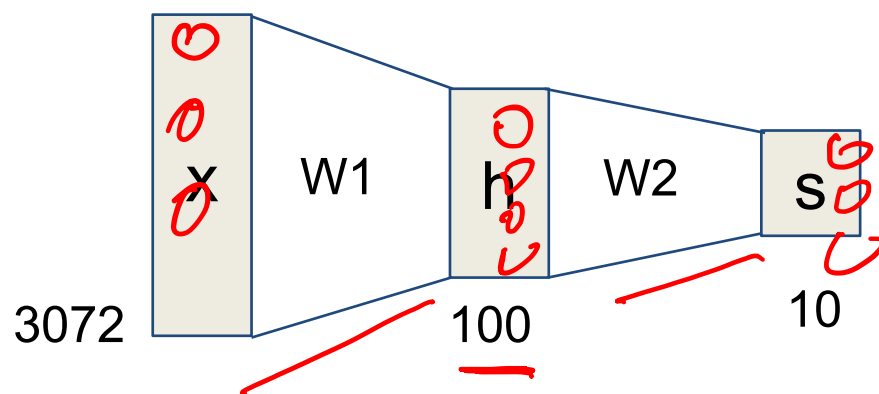
Handwritten mathematical expressions in red ink:

$$\mathbb{R}^{100} \Rightarrow h = \left[ \begin{array}{c} \vdots \\ \vdots \end{array} \right]_{100 \times 3072}$$
$$s = W_2 \left[ \begin{array}{c} \vdots \\ \vdots \end{array} \right]_{10 \times 100}$$

# Neural networks: without the brain stuff

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

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



# Neural networks: without the brain stuff

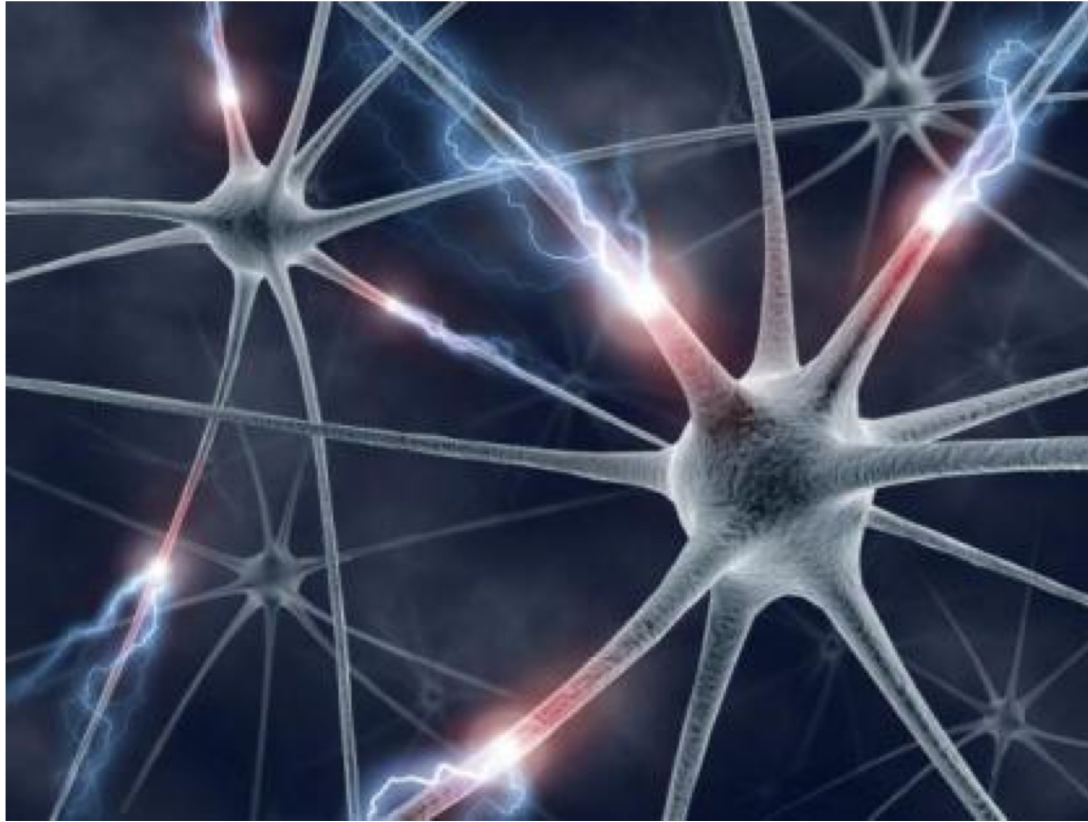
(**Before**) Linear score function:  $f = Wx$

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

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

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

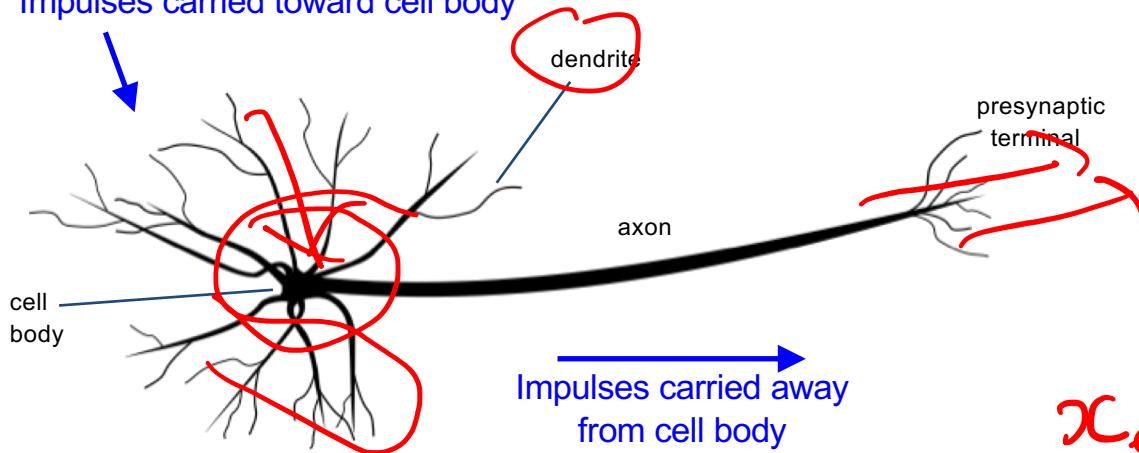
```
1 import numpy as np
2 from numpy.random import randn
3
4 N, D_in, H, D_out = 64, 1000, 100, 10
5 x, y = randn(N, D_in), randn(N, D_out)
6 w1, w2 = randn(D_in, H), randn(H, D_out)
7
8 for t in range(2000):
9     h = 1 / (1 + np.exp(-x.dot(w1)))
10    y_pred = h.dot(w2)
11    loss = np.square(y_pred - 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
```



[This image](#) by [Fotis Bobolas](#) is licensed under [CC-BY 2.0](#)

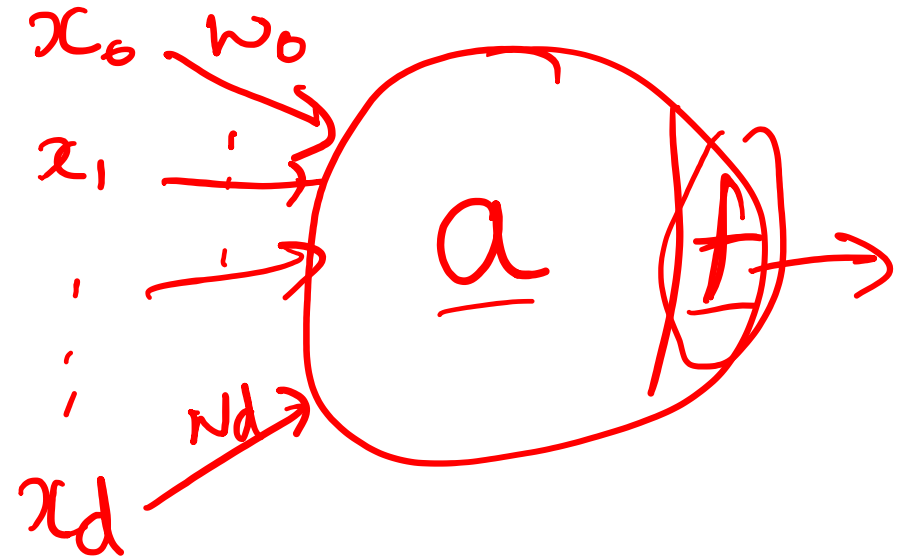
Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

Impulses carried toward cell body



This image by Felipe Perucho is licensed under [CC-BY 3.0](https://creativecommons.org/licenses/by/3.0/)

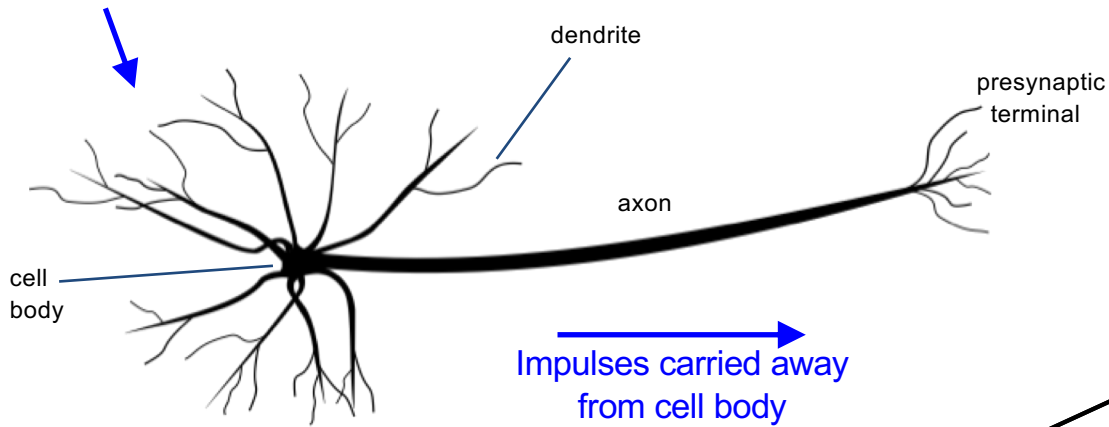
$$a = \sum_j w_j x_j$$
$$= \underline{\vec{w}}^T \underline{x}$$



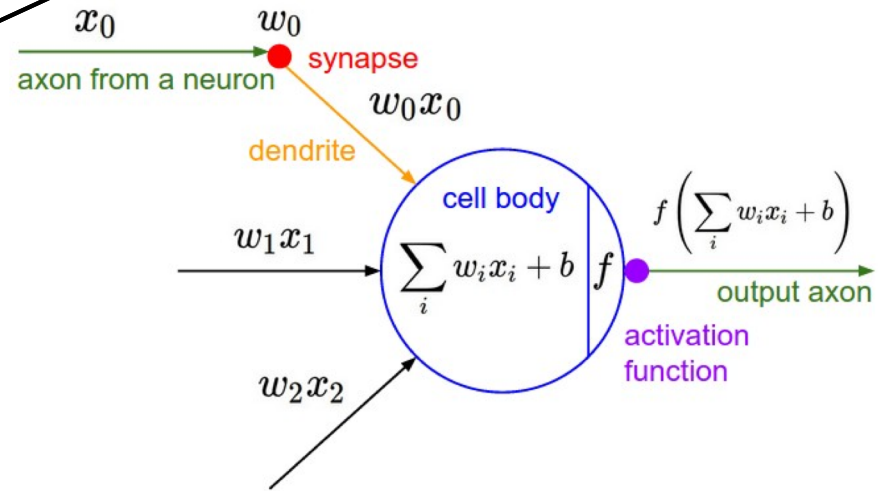
$$y = f(a)$$



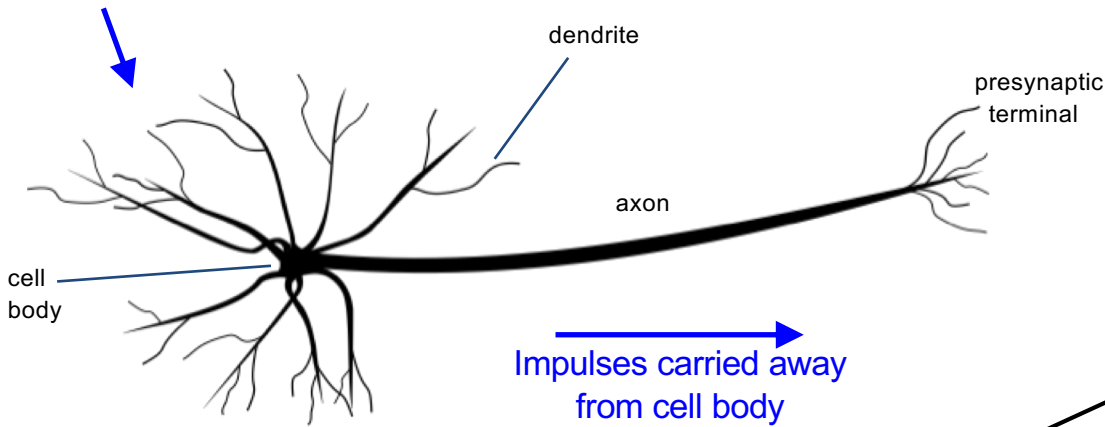
Impulses carried toward cell body



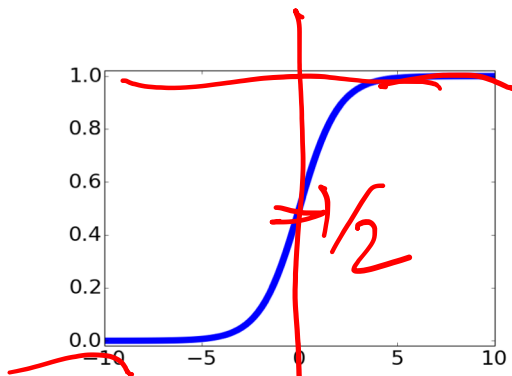
[This image](#) by Felipe Perucho is licensed under [CC-BY 3.0](#)



Impulses carried toward cell body



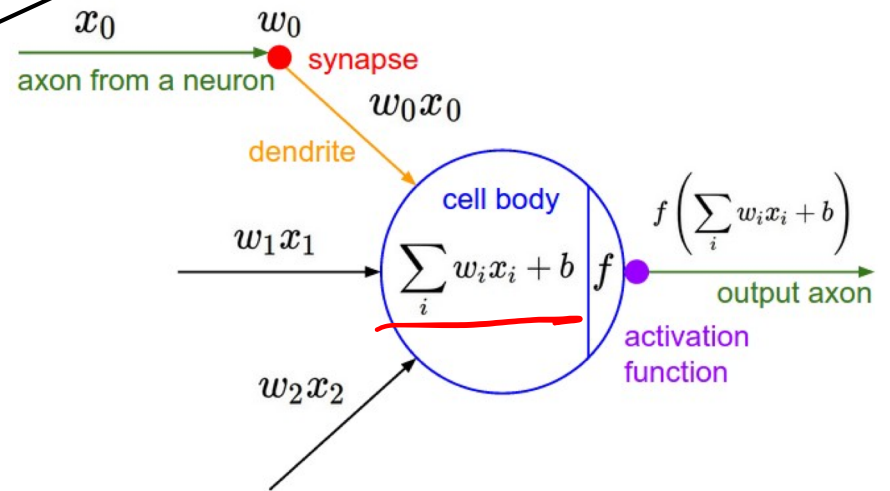
This image by Felipe Perucho is licensed under [CC-BY 3.0](https://creativecommons.org/licenses/by/3.0/)



sigmoid activation function

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

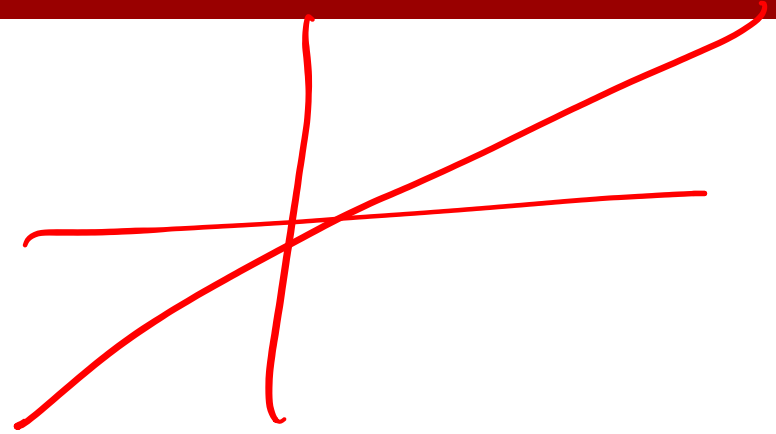
$$f(a) = \frac{1}{1 + e^{-a}}$$



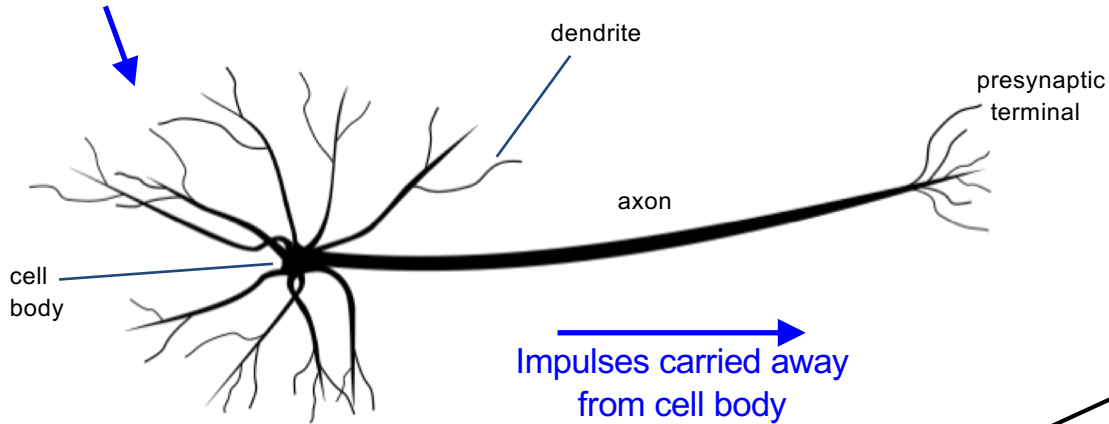
$$f(a) = a$$

$$f(y) = \vec{w}^T \vec{x}$$

$$f(a) = \max(0, a)$$



Impulses carried toward cell body



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

```
class Neuron:
```

```
# ...
```

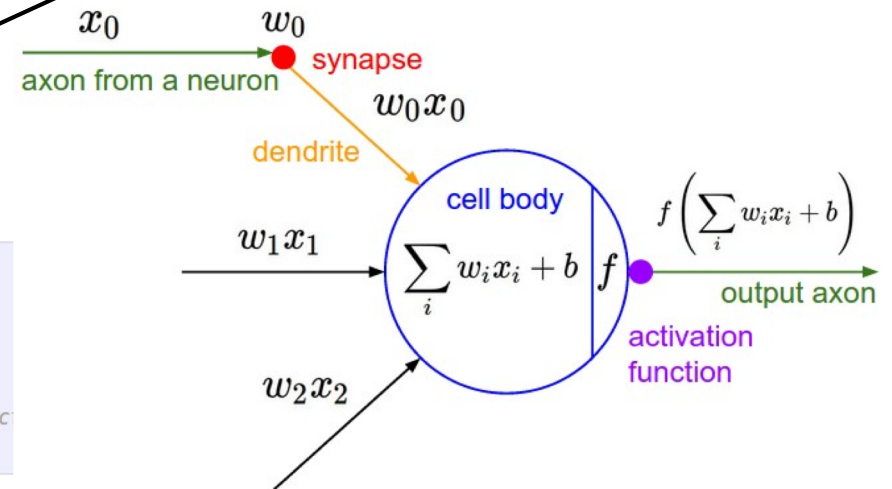
```
def neuron_tick(inputs):
```

```
    """ 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
```



# Be very careful with your brain analogies!

## **Biological Neurons:**

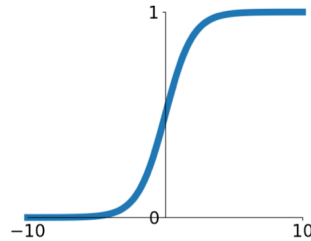
- Many different types
- Dendrites can perform complex non-linear computations
- Synapses are not a single weight but a complex non-linear dynamical system
- Rate code may not be adequate

[Dendritic Computation. London and Hausser]

# Activation functions

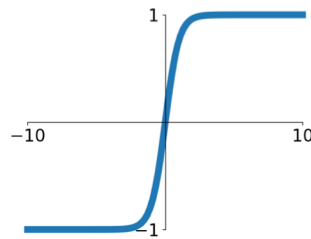
**Sigmoid**

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



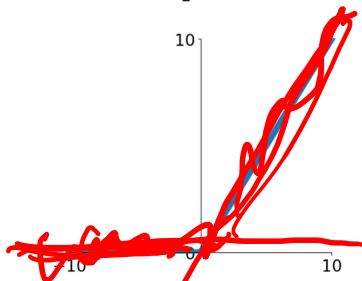
**tanh**

$$\tanh(x)$$



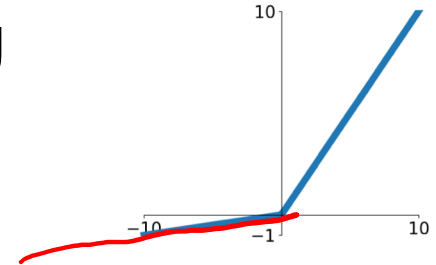
**ReLU**

$$\max(0, x)$$



**Leaky ReLU**

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

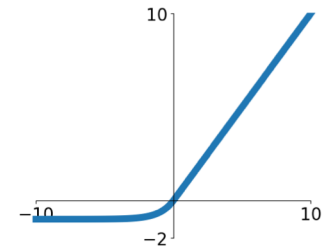


**Maxout**

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

**ELU**

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

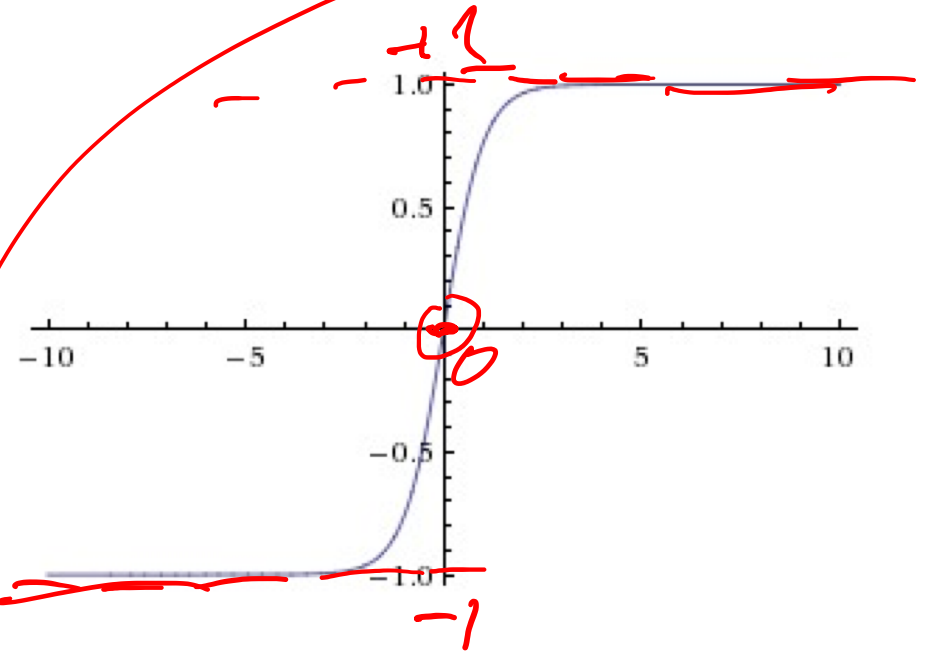
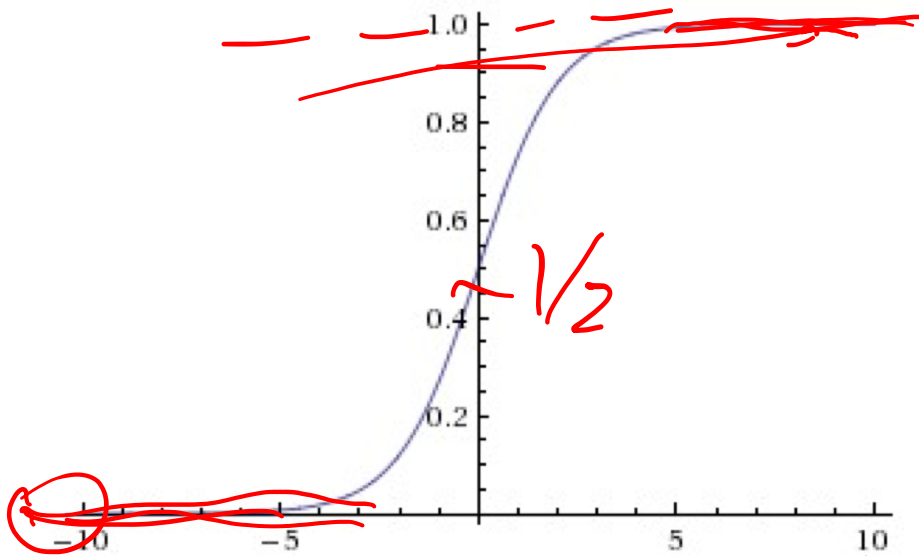


# Activation Functions

- sigmoid vs tanh

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

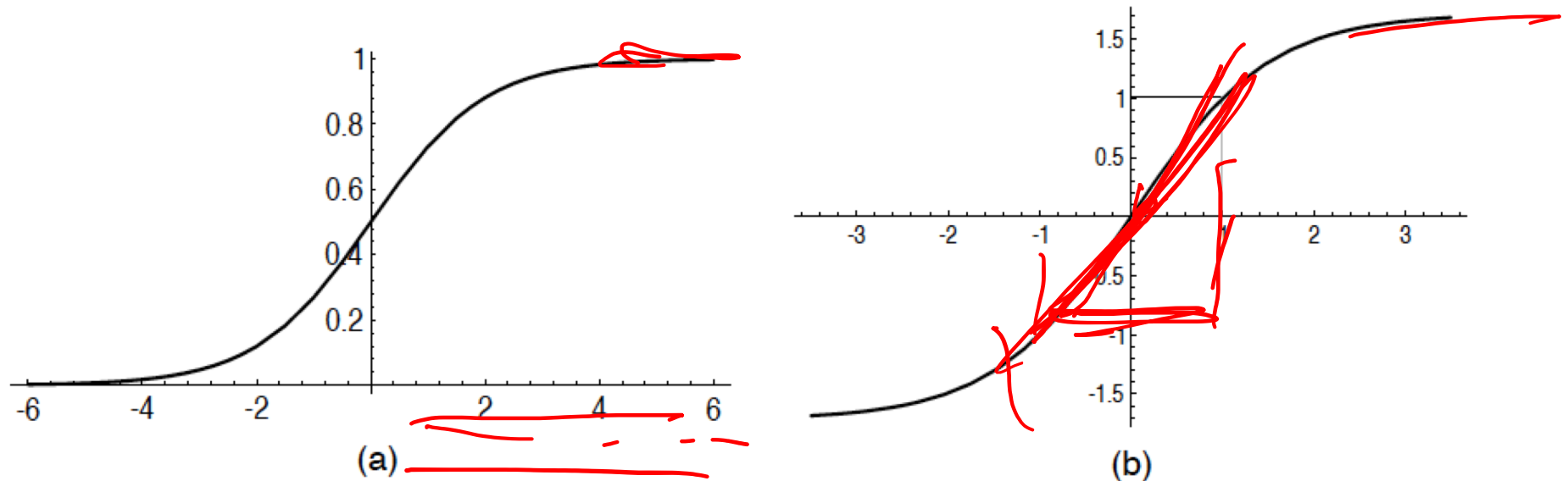
$$\tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$$



$$\sigma(\sigma(\sigma(a)))$$

$$\Rightarrow = \underline{\underline{2\sigma(2a) - 1}}$$

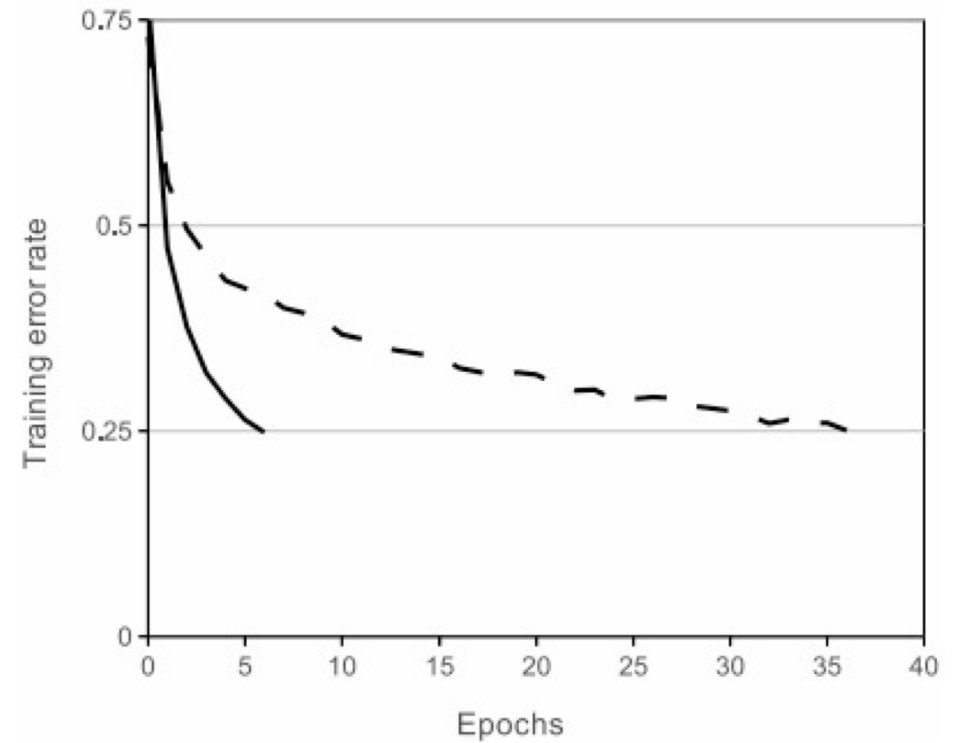
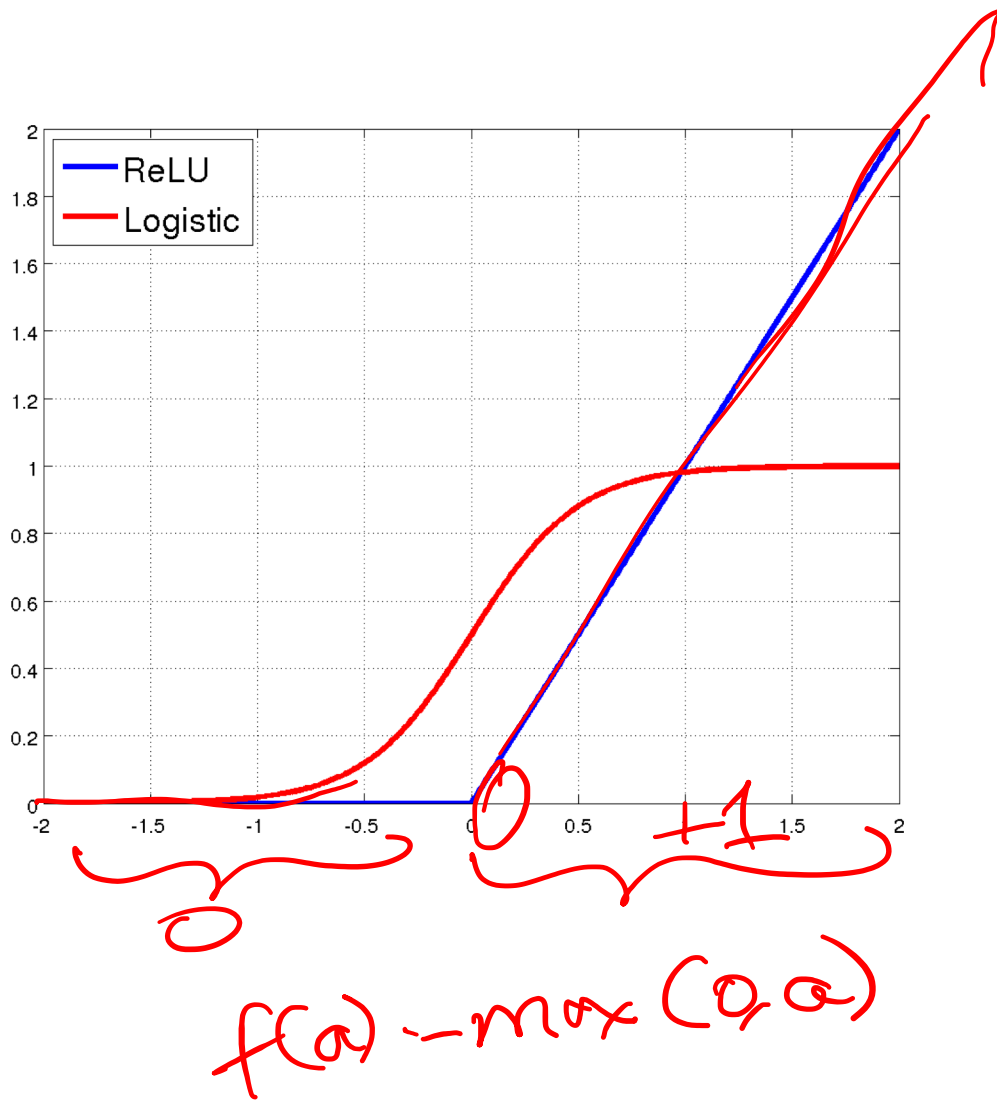
# A quick note



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



# Rectified Linear Units (ReLU)



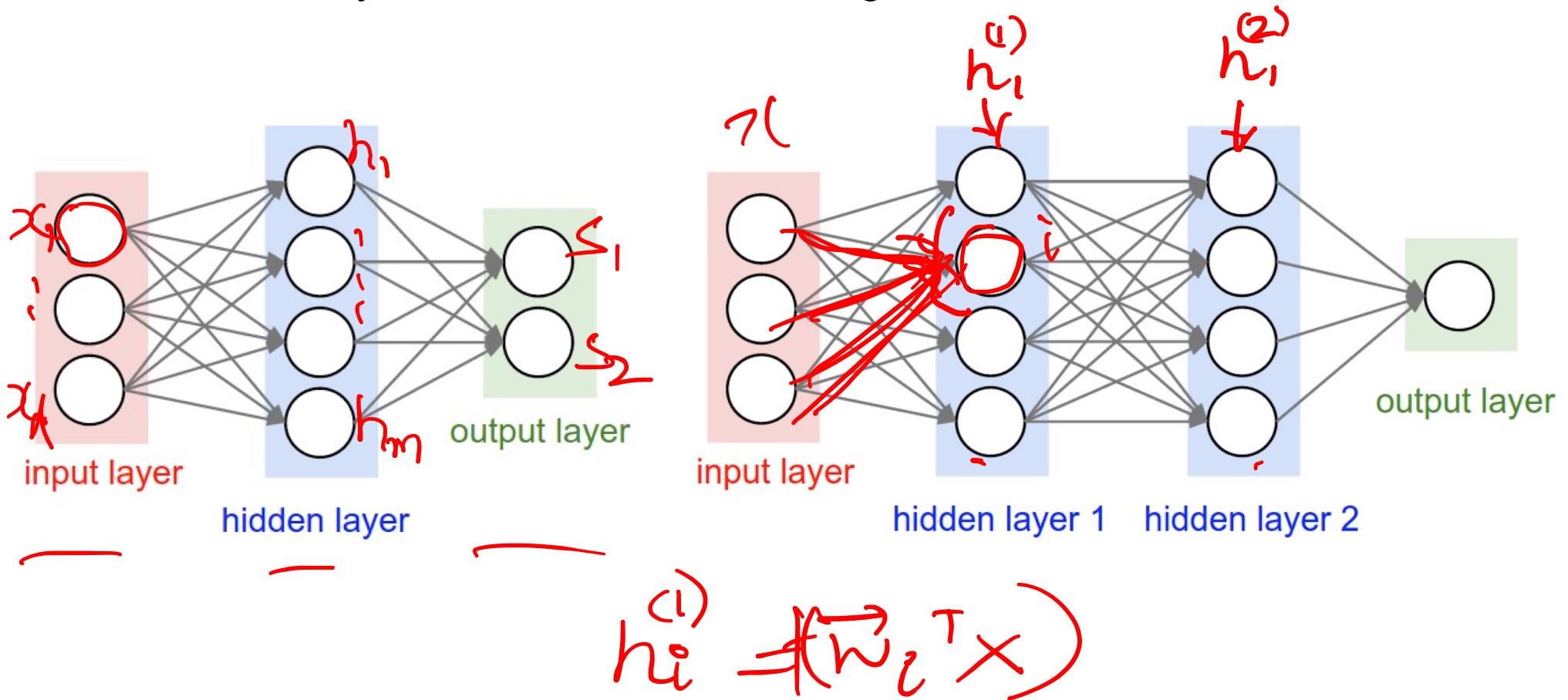
[Krizhevsky et al., NIPS12]

# Limitation

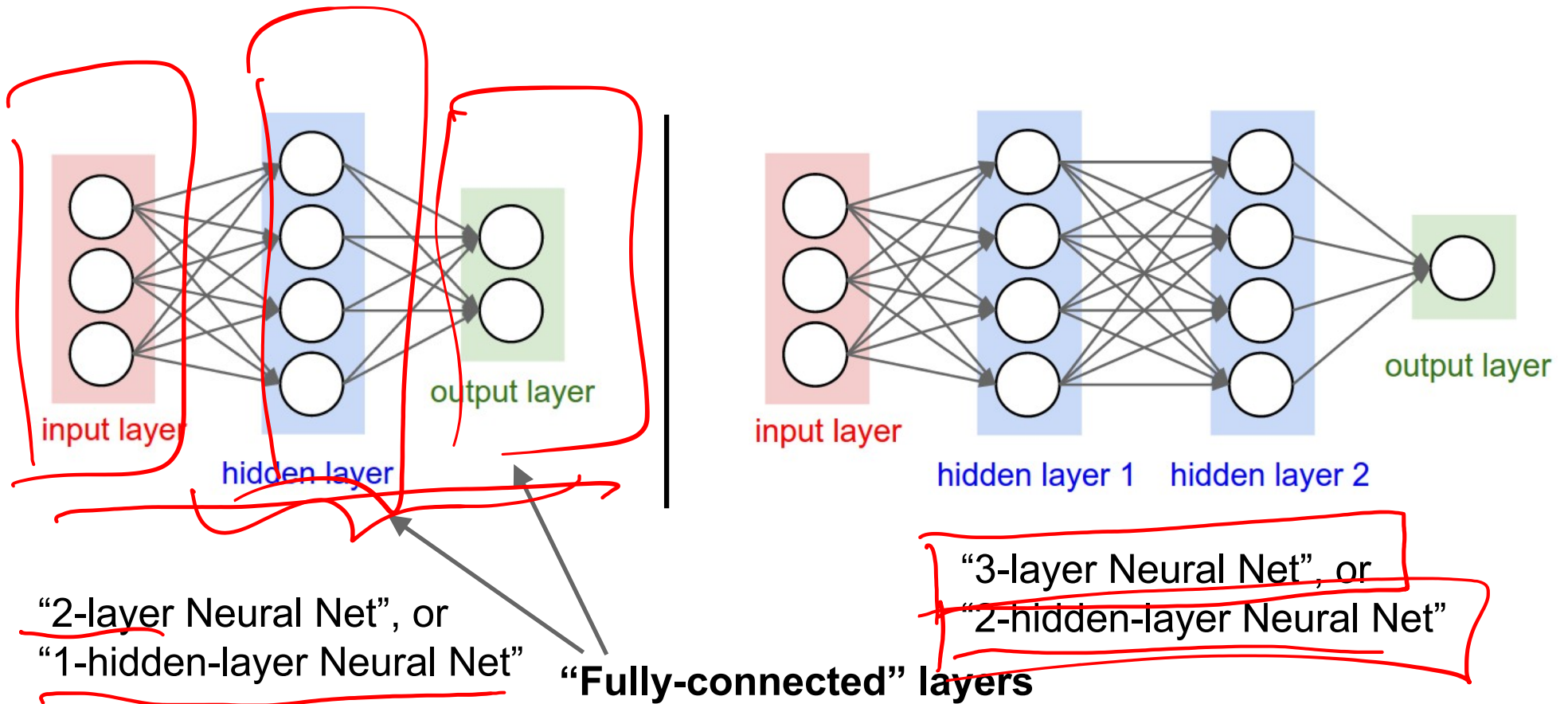
- A single “neuron” is still a linear decision boundary
- What to do?
- Idea: Stack a bunch of them together!

# Multilayer Networks

- Cascade Neurons together
- The output from one layer is the input to the next
- Each Layer has its own sets of weights



# Neural networks: Architectures



# Demo Time

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



# Optimization

$$\min_w L(w; D)$$

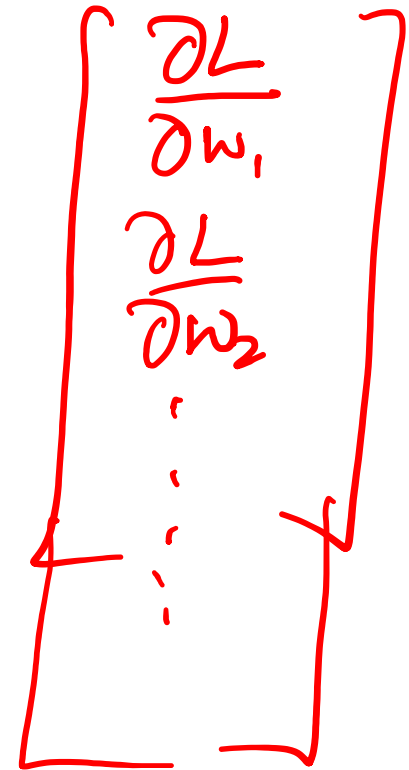
Strategy: **Follow the slope**



## Strategy: **Follow the slope**

In 1-dimension, the derivative of a function:

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



A handwritten red vector notation representing a gradient vector. It consists of a large square bracket on the right side, enclosing a vertical list of terms. From top to bottom, the terms are:  $\frac{\partial L}{\partial w_1}$ ,  $\frac{\partial L}{\partial w_2}$ , and a vertical ellipsis ( $\vdots$ ).

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

```
# Vanilla Gradient Descent
```

```
while True:
```

```
    [weights_grad = evaluate_gradient(loss_fun, data, weights)] ← backprop  
    weights += - step_size * weights_grad # perform parameter update
```

$$W^{(0)} = \text{init}$$

for  $t=1 \dots \text{times}$

$$\vec{w}^{(t+1)} = \vec{w}^t - \eta \boxed{\nabla_w L}$$

