

# CS 4803 / 7643: Deep Learning

Topics:

- Regularization
- Neural Networks

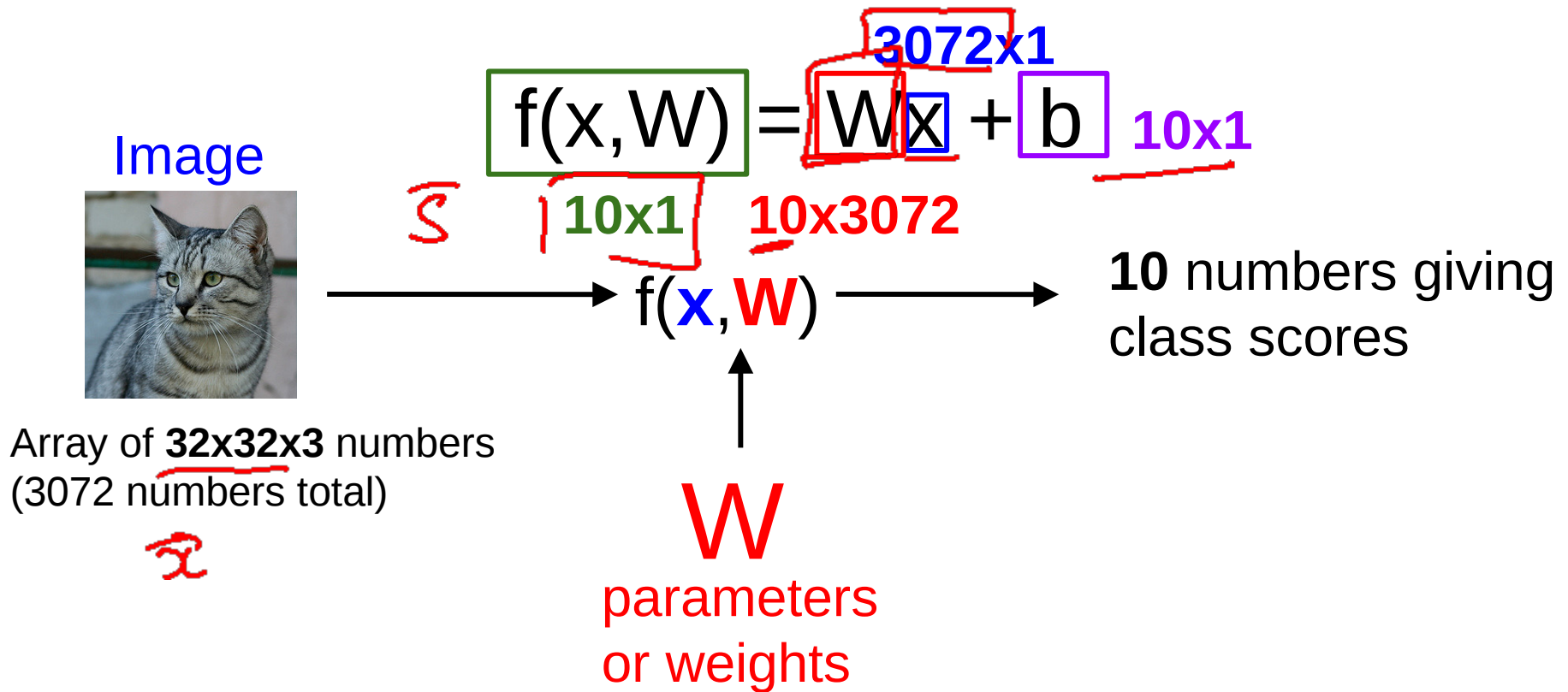
Dhruv Batra  
Georgia Tech

# Administrativa

- PS1/HW1 out
  - Due: 09/09, 11:59pm
  - Asks about topics being covered now
  - Caveat: one more (extra credit) problem to be added
  - Please please please please start early /
  - More details next class

# Recap from last time

# Parametric Approach: Linear Classifier



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

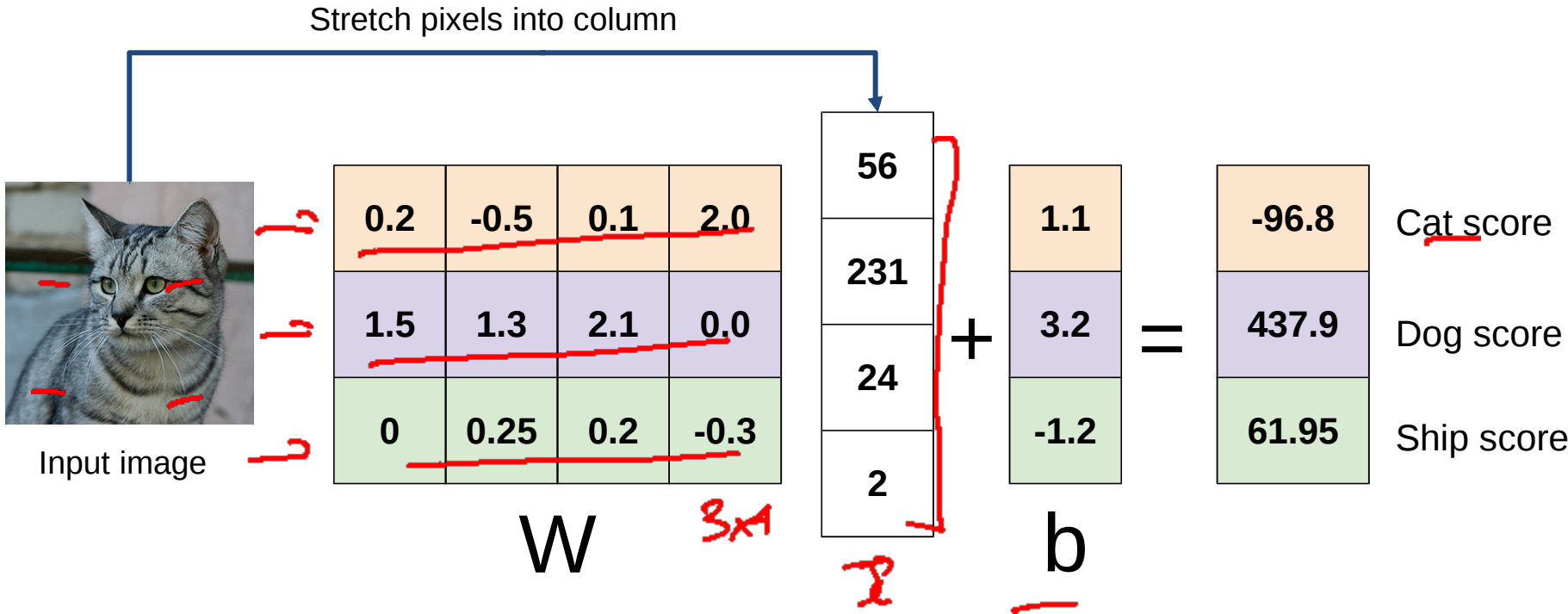
Stretch pixels into column



Input image

56
231
24
2

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



0.2	-0.5	0.1	2.0
1.5	1.3	2.1	0.0
0	0.25	0.2	-0.3

$W$

56
231
24
2

$x_i$

1.1
3.2
-1.2

$b$



0.2	-0.5	0.1	2.0	1.1
1.5	1.3	2.1	0.0	3.2
0	0.25	0.2	-0.3	-1.2

$W$        $b$

new, single  $W$

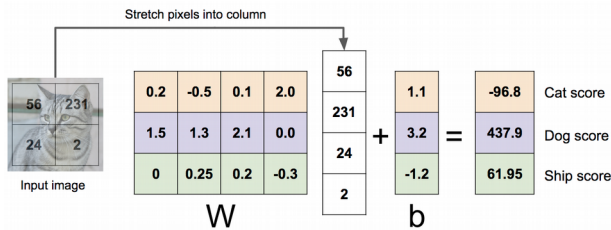
56
231
24
2
1

$x_i$

# Linear Classifier: Three Viewpoints

## Algebraic Viewpoint

$$f(x, W) = \underline{Wx} + \underline{b}$$



## Visual Viewpoint

One template per class



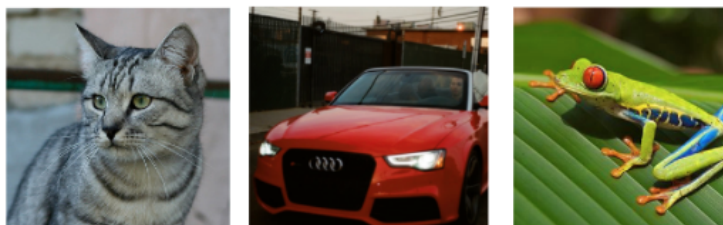
## Geometric Viewpoint

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

## TODO:

1. Define a loss function that quantifies our unhappiness with the scores across the training data.
2. Come up with a way of efficiently finding the parameters that minimize the loss function. (**optimization**)

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

# Softmax vs. SVM

$P(y_i | x_i, \bar{w})$

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

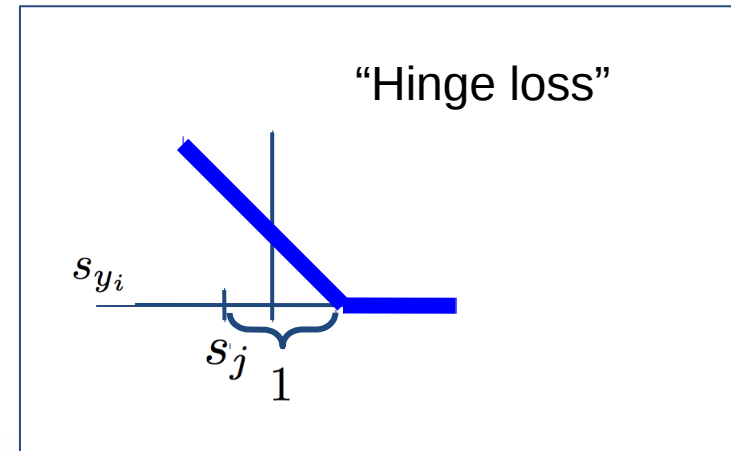
$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

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



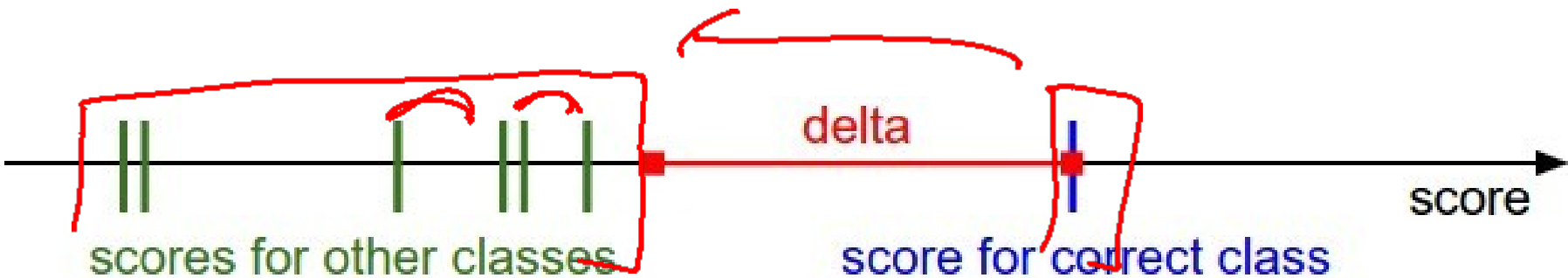
cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>

### Multiclass SVM loss:



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

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



## 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	<u>-1.7</u>

# 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  $\geq 0$

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**

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

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

Probabilities  
must be  $\geq 0$

Probabilities  
must sum to 1

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

exp

<b>24.5</b>
<b>164.0</b>
<b>0.18</b>

normalize

<b>0.13</b>
<b>0.87</b>
<b>0.00</b>

unnormalized  
probabilities

probabilities

# Softmax Classifier (Multinomial Logistic Regression)



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

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

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

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

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

# 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

# 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

cat	3.2
car	5.1
frog	-1.7

Unnormalized log-probabilities / logits

exp

cat	24.5
car	164.0
frog	0.18

unnormalized probabilities

normalize

cat	0.13
car	0.87
frog	0.00

probabilities

compare  
min  
KL

cat	1.00
car	0.00
frog	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

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**

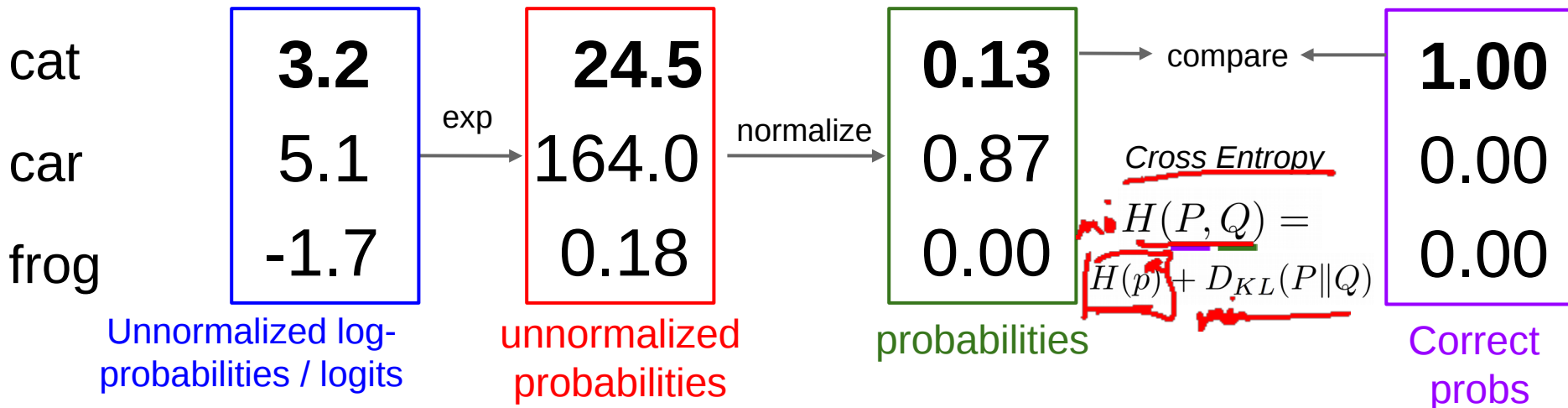
$$s = f(x_i; W)$$

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

Probabilities must be  $\geq 0$

Probabilities must sum to 1

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



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

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


IID  $\sim P^x$

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

# Plan for Today

- Regularization
- Neural Networks

# Regularization

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


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



# Regularization

*objective*

$$L(W) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)}_{\substack{-\log P(D|\vec{w}) \\ \text{"model fit"}}} + \underbrace{\lambda R(W)}_{-\log P(\vec{w})}$$

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

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

# Regularization

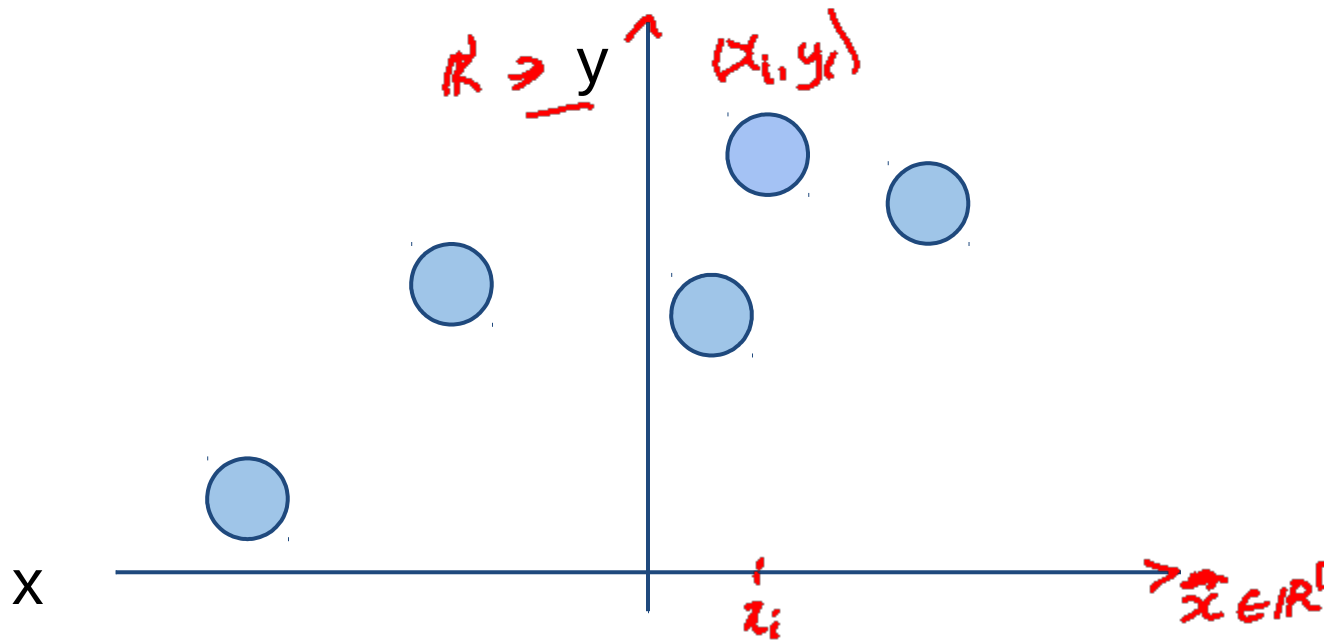
$\lambda$  = regularization strength  
(hyperparameter)

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

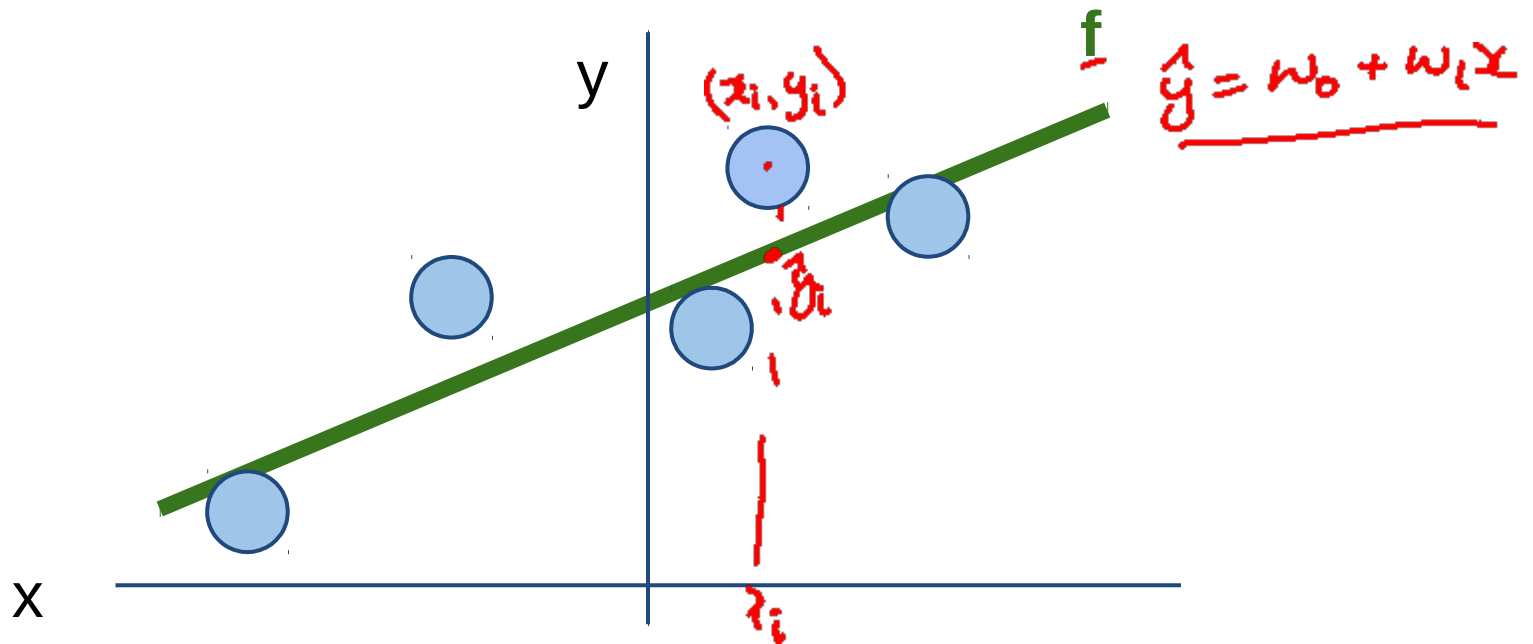
**Data loss:** Model predictions should match training data

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

# Regularization Intuition in Polynomial Regression



# Polynomial Regression



# Polynomial Regression

$$\hat{y} = \underset{\downarrow}{w_0} + \underset{\downarrow}{w_1} x \leftarrow \boxed{1^{\text{st}} \text{ order poly}}$$

$$= w_0 + w_1 x + w_2 x^2$$

$$= \underset{\downarrow}{w_0} + \dots + \underset{\uparrow}{w_d x^d} \leftarrow \boxed{d^{\text{th}} \text{ poly}}$$

$$= [w_0 \dots w_d] \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^d \end{bmatrix}$$

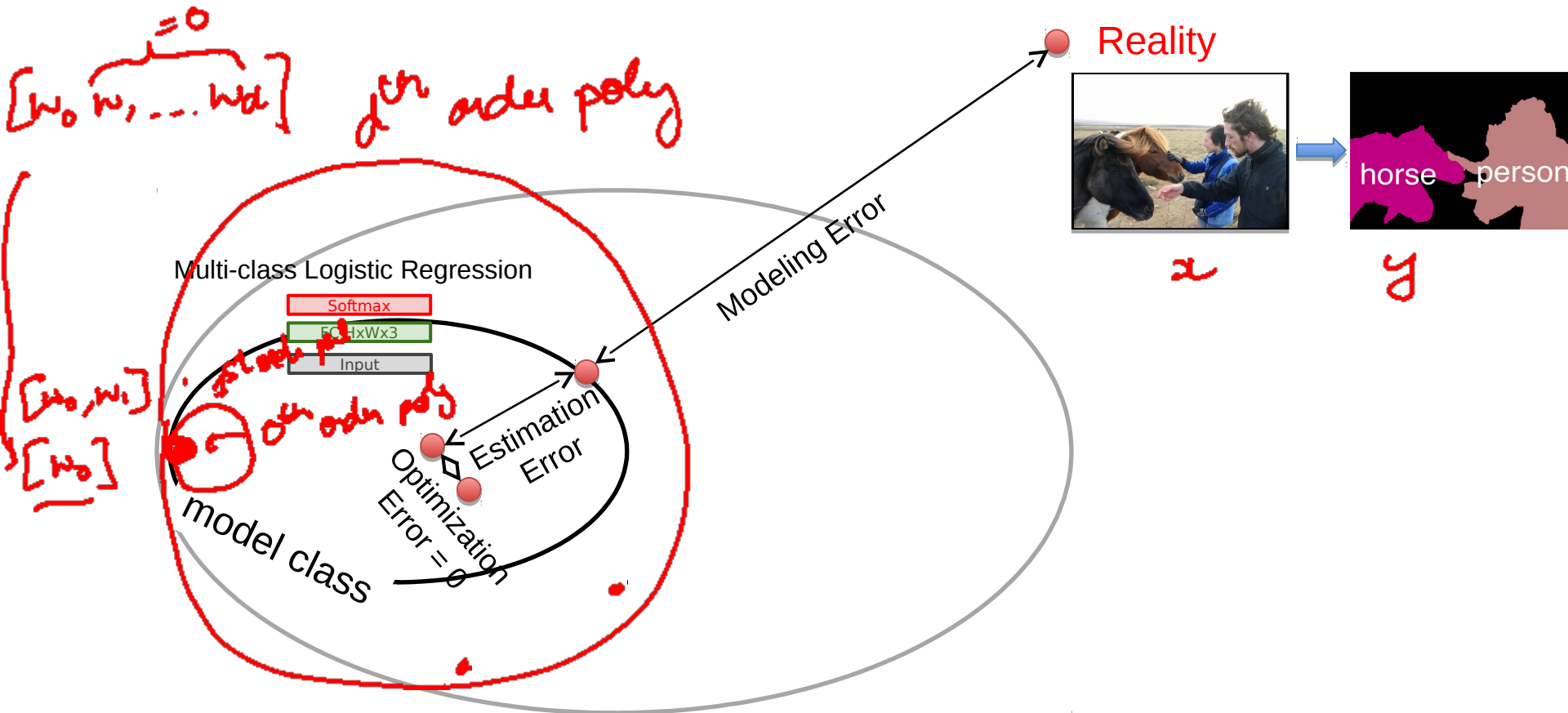
$$= \underbrace{\vec{w}}^T \underbrace{\vec{\phi}(x)}$$

linear in  $\vec{w}$

$$\min_{\vec{w}} L(w, \delta) = \frac{1}{N} \sum_{i=1}^N (y_i - \vec{w}^T \vec{\phi}(x_i))^2$$

# Polynomial Regression

# Error Decomposition



# Polynomial Regression

- Demo: <https://arachnoid.com/polysolve/>
- You are a scientist studying runners.
  - You measure average speeds of the best runners at different ages.
- Data: Age (years), Speed (mph)

–	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

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

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

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

# 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

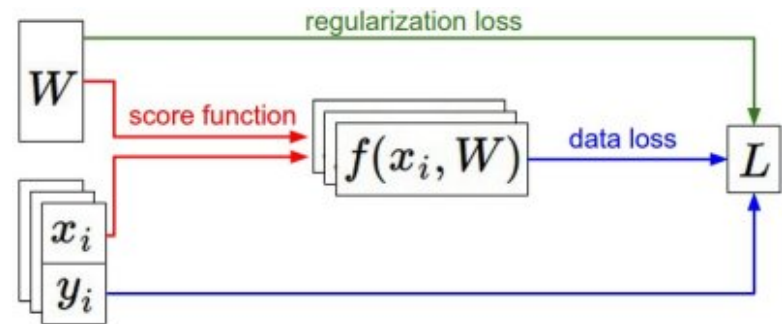
# Recap

- 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**:

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

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

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



# Recap

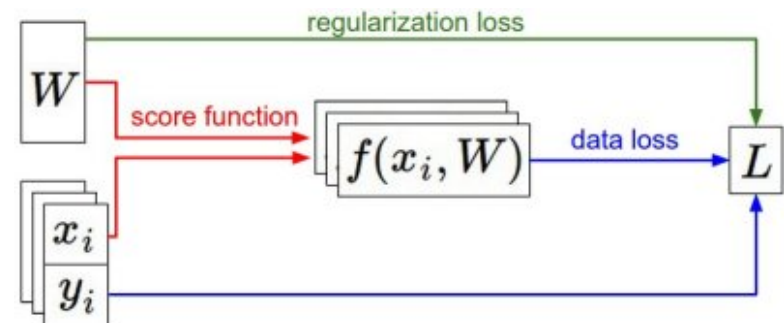
How do we find the best  $W$ ?

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

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

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

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



Next: Neural Networks

# So far: Linear Classifiers



"Row Input"  
x



$$f(x) = \underline{W}x$$

Class  
scores





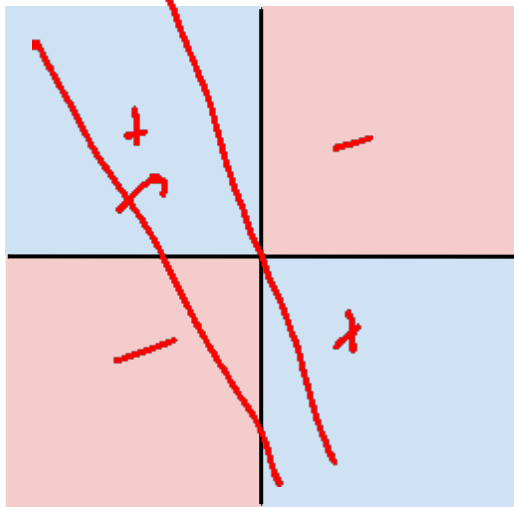
# Hard cases for a linear classifier

**Class 1:**

First and third quadrants

**Class 2:**

Second and fourth quadrants

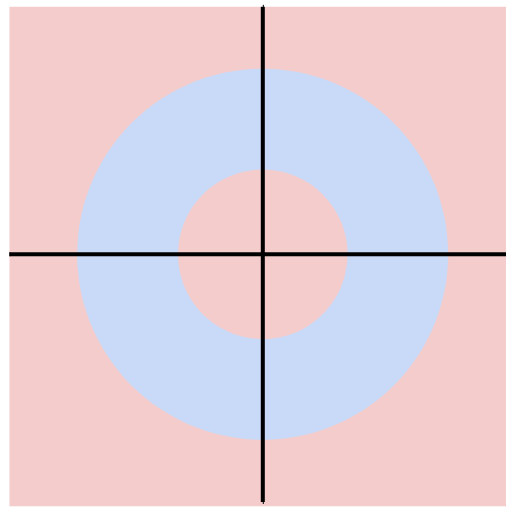


**Class 1:**

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

**Class 2:**

Everything else

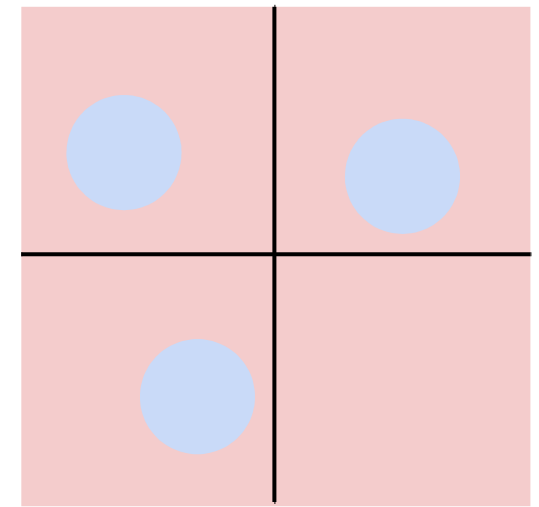


**Class 1:**

Three modes

**Class 2:**

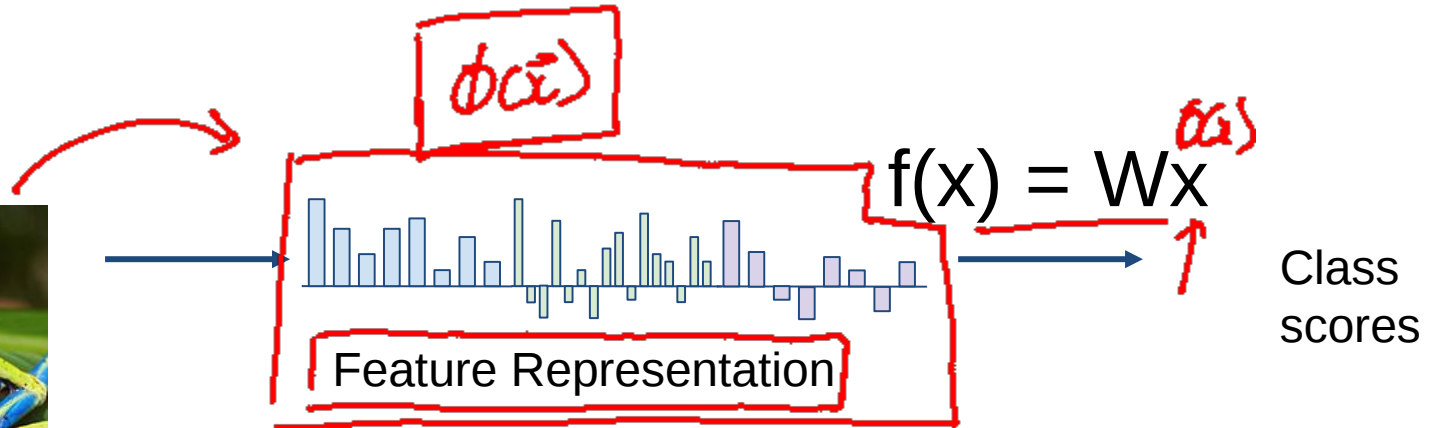
Everything else



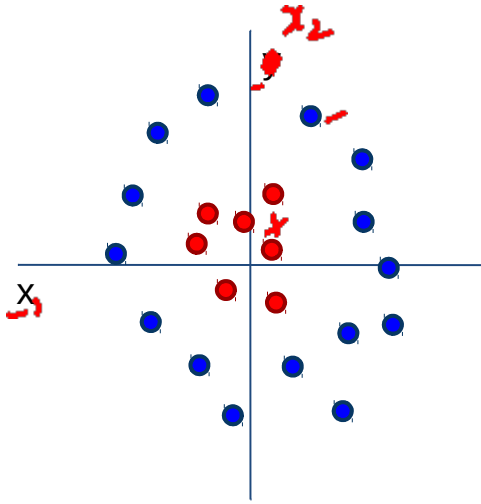
# Aside: Image Features



Raw Input  
 $\vec{x}$

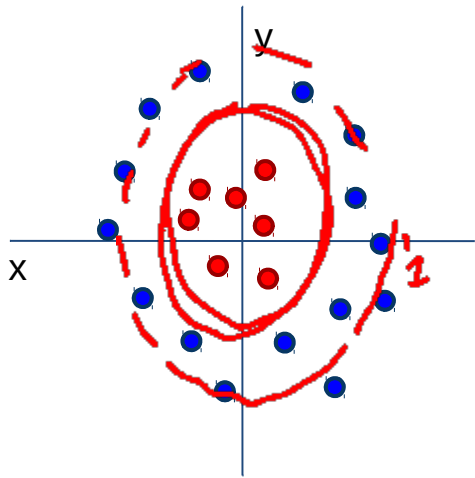


# Image Features: Motivation



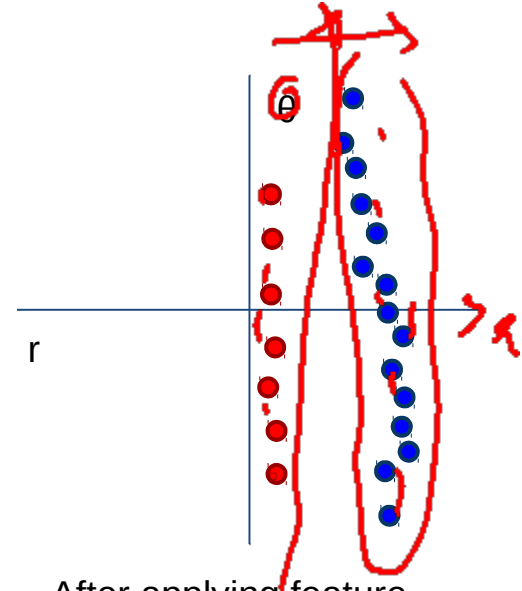
Cannot separate red  
and blue points with  
linear classifier

# Image Features: Motivation



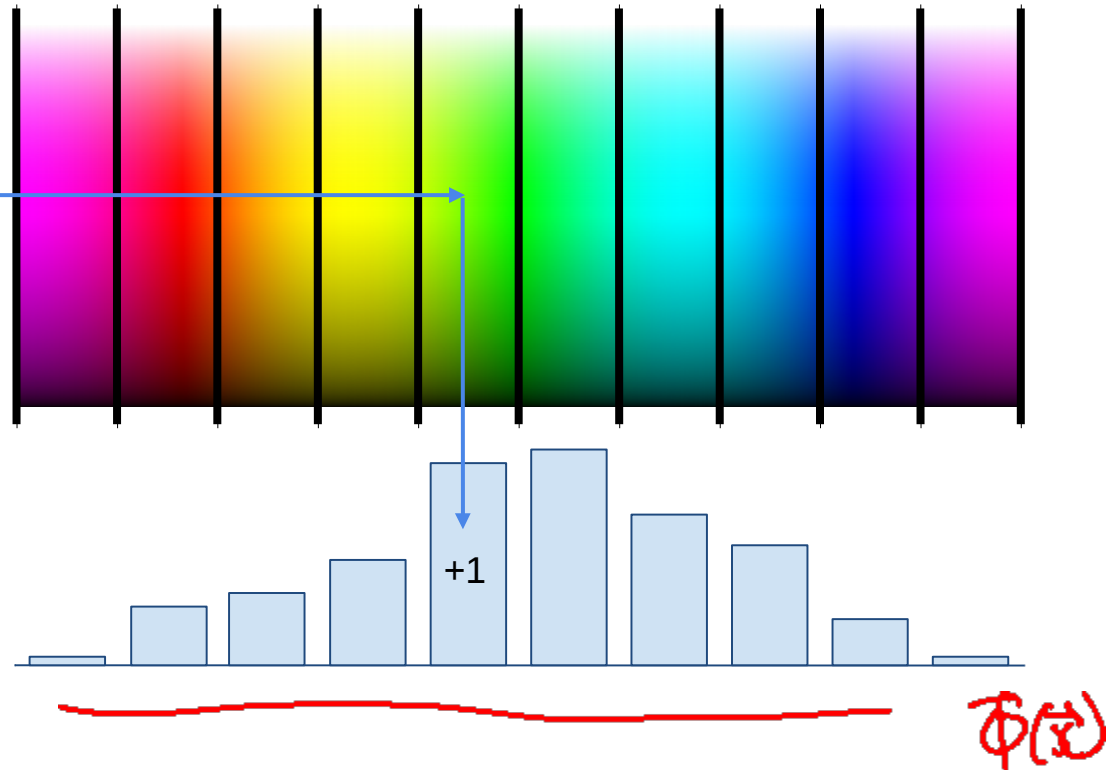
Cannot separate red and blue points with linear classifier

$$f(x, y) = (\underline{r(x, y)}, \underline{\theta(x, y)})$$

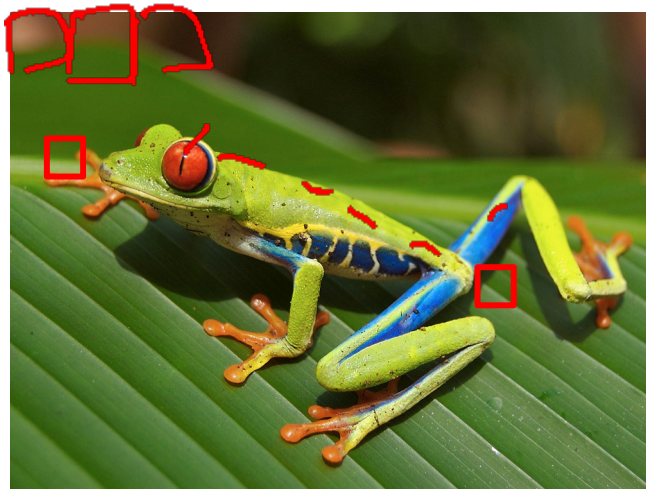


After applying feature transform, points can be separated by linear classifier

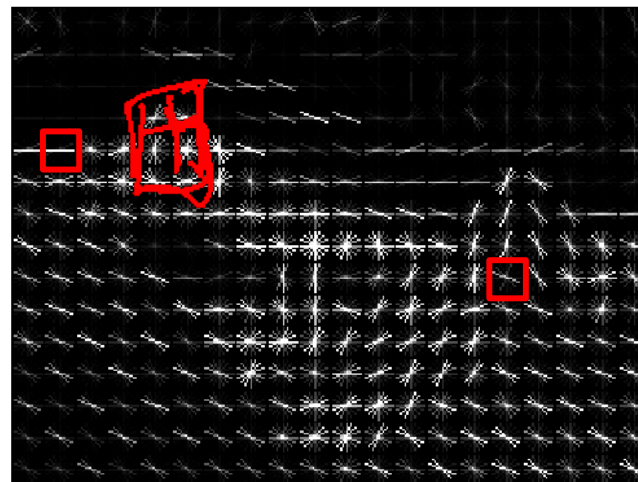
# Example: Color Histogram



# Example: Histogram of Oriented Gradients (HoG)



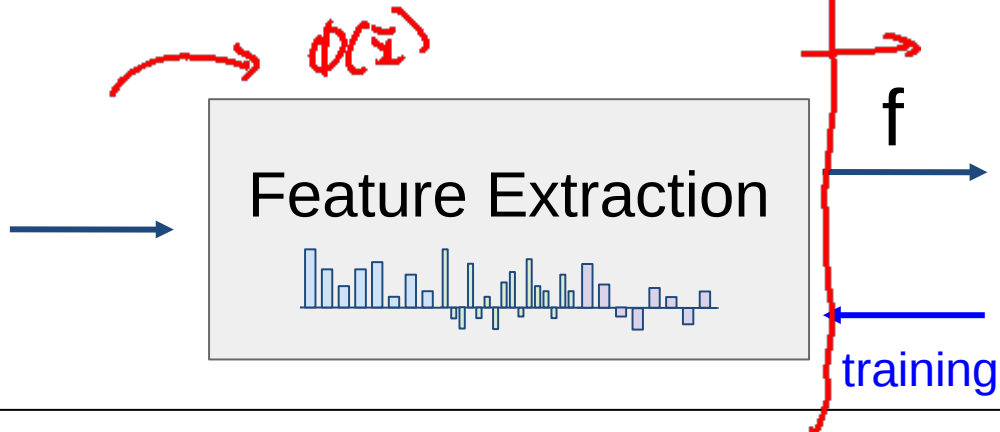
Divide image into 8x8 pixel regions  
Within each region quantize edge  
direction into 9 bins



Example: 320x240 image gets divided  
into 40x30 bins; in each bin there are  
9 numbers so feature vector has  
 $30 \times 40 \times 9 = 10,800$  numbers

Lowe, "Object recognition from local scale-invariant features", ICCV 1999  
Dalal and Triggs, "Histograms of oriented gradients for human detection," CVPR 2005

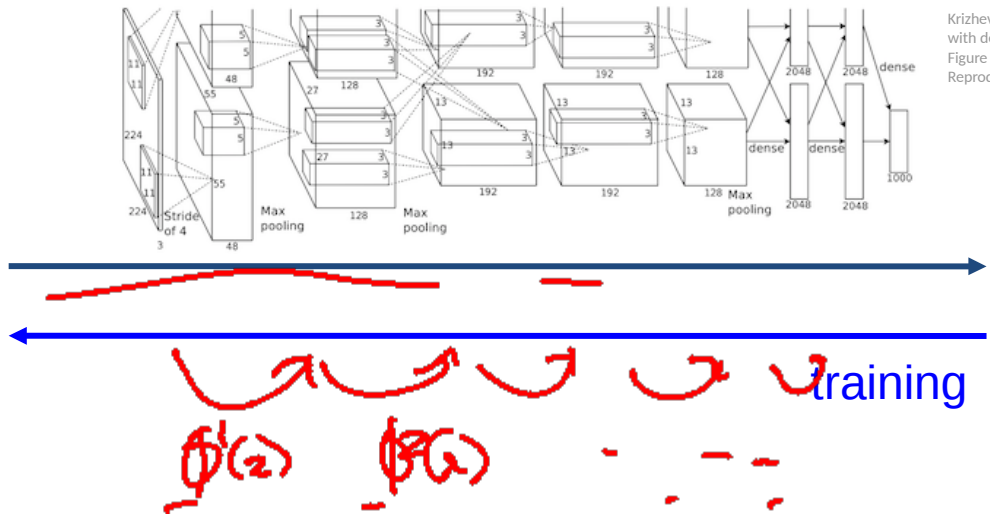
# Image features vs Neural Nets



10 numbers giving scores for classes



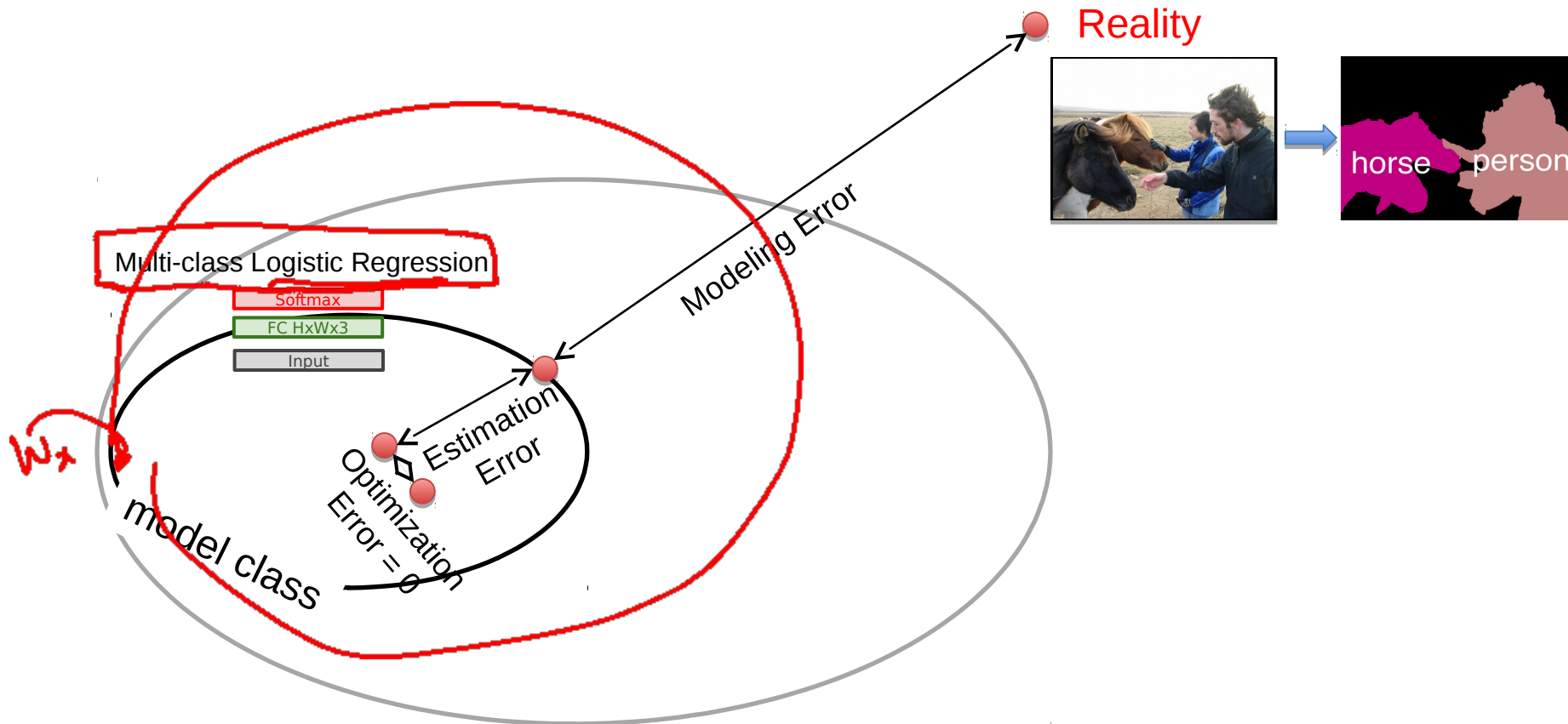
$x$



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

# Error Decomposition





# Neural networks: without the brain stuff

$$H = \{h: X \rightarrow Y\}$$

(Before) Linear score function:

$$f = \underline{W}x$$

The diagram shows a red box containing  $w_2$  and  $w_1$  with a plus sign between them, and an  $x$  to the right. Below this, another red box contains  $w_3$  and an  $x$  to the right. A dashed line is on the left, and a red arrow points from the right towards the diagram.

# Neural networks: without the brain stuff



(Before) Linear score function:

(Now) 2-layer Neural Network

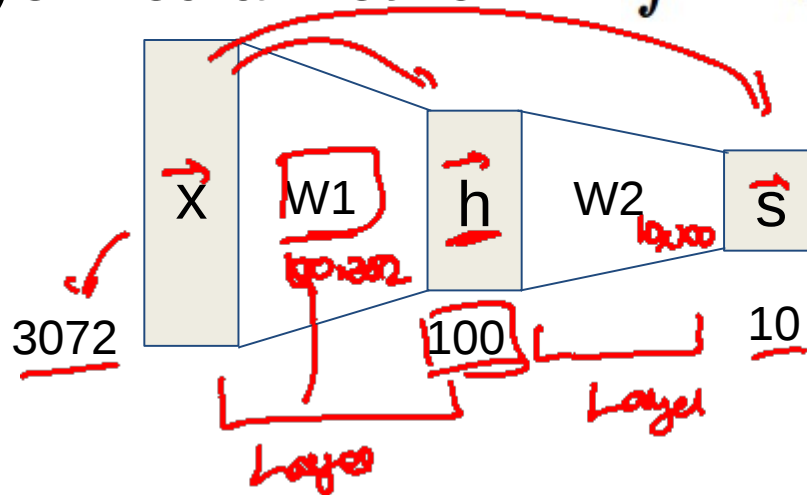
$$f = Wx$$
$$f = W_2 \max(0, W_1 \vec{x})$$

Diagram illustrating the structure of a 2-layer neural network. A vector  $\vec{x}$  is multiplied by a weight matrix  $W_1$  to produce a hidden layer output. This output is then passed through a  $\max(0, \cdot)$  activation function. The result is multiplied by a weight matrix  $W_2$  to produce the final output  $f$ . The diagram shows a box containing a vertical list of nodes:  $\text{node } h_1$ ,  $i$ , and  $\text{node } h_2$ . A weight  $w_3$  is shown below the box, and a weight  $w_2$  is shown to the left of the box. Arrows indicate the flow of information from  $\vec{x}$  through  $W_1$  to the hidden layer, through the activation function, and through  $W_2$  to the output  $f$ .

# Neural networks: without the brain stuff

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

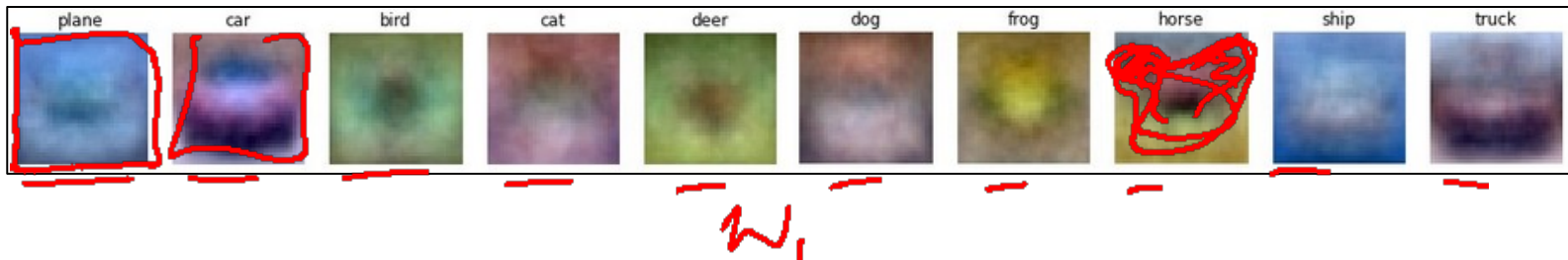
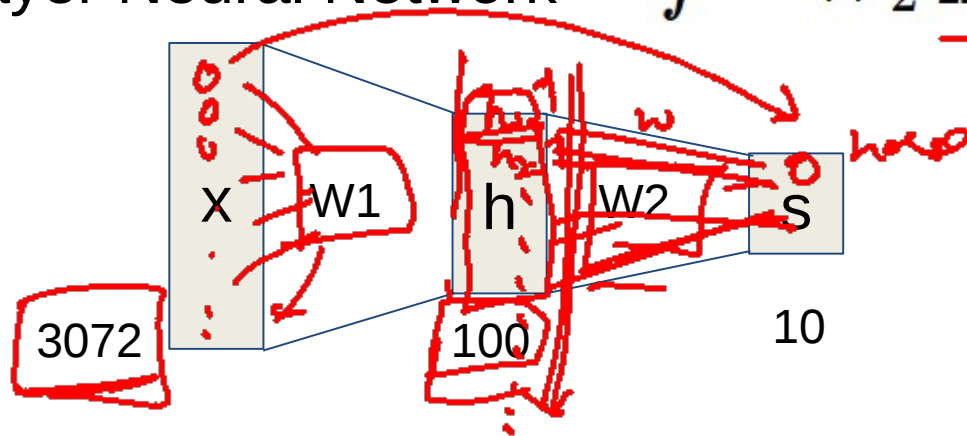
(**Now**) 2-layer Neural Network  $f = W_2 \max(0, W_1 x)$



# 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

$$L(\{w_1, \dots, w_3\}, D)$$

$$L_i(\underline{g(x_i)}, \underline{y_i})$$

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

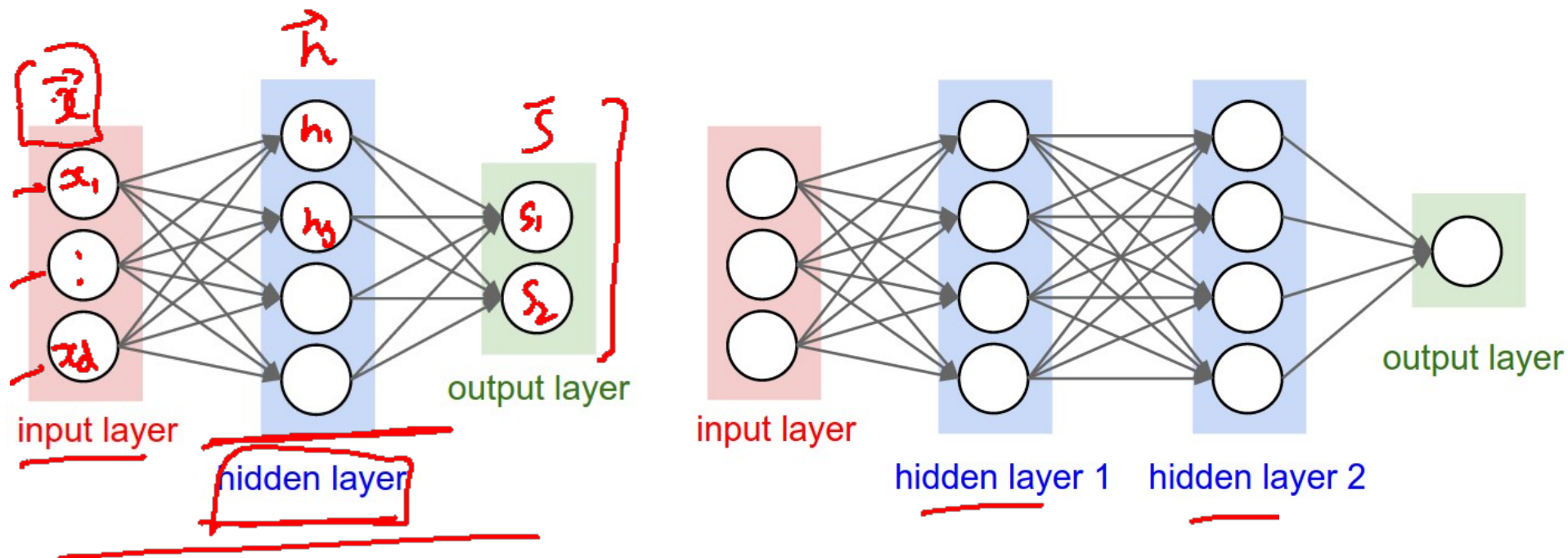
(**Now**) 2-layer Neural Network  
or 3-layer Neural Network

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

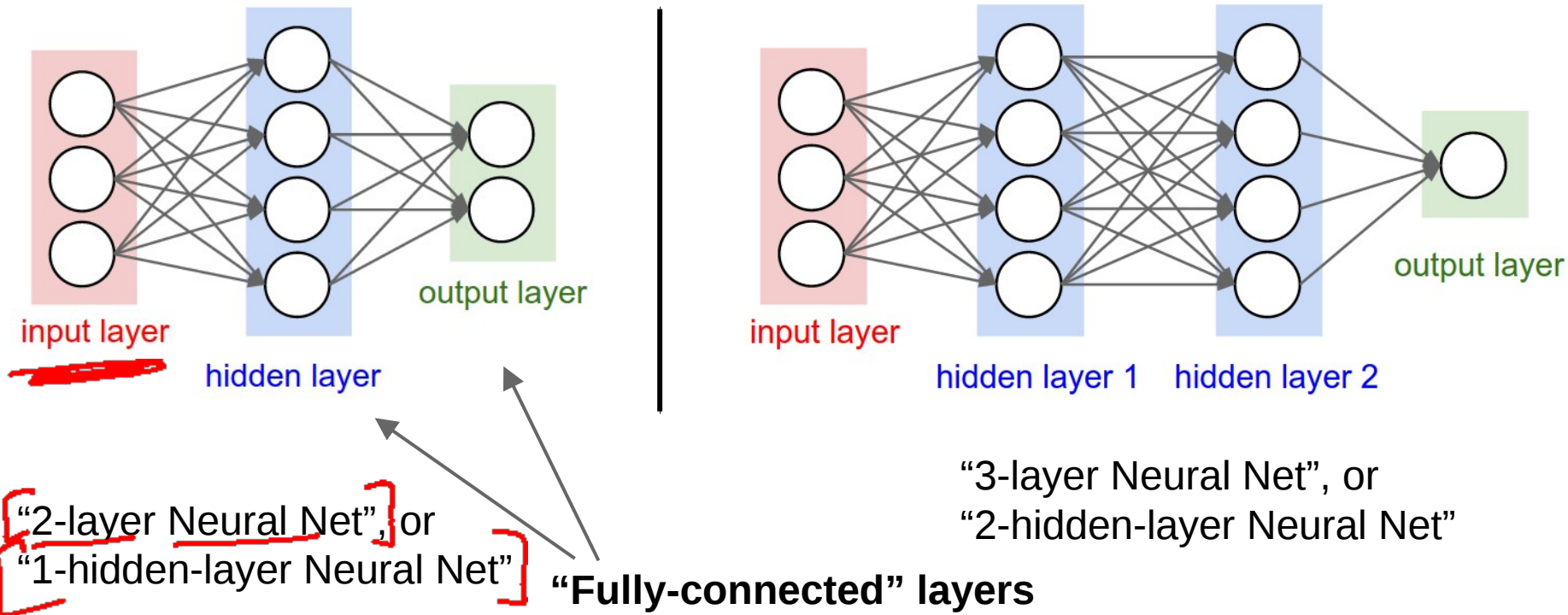
$$\vec{x} \xrightarrow{w_1} \vec{h}^{(1)} \xrightarrow{w_2} \vec{h}^{(2)} \xrightarrow{w_3} \vec{s}$$

# Multilayer Networks

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



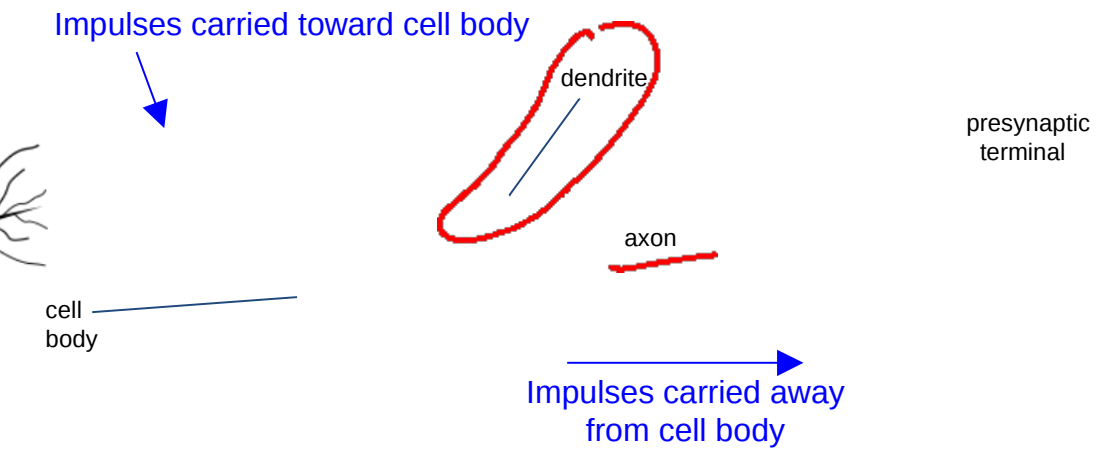
# Neural networks: Architectures



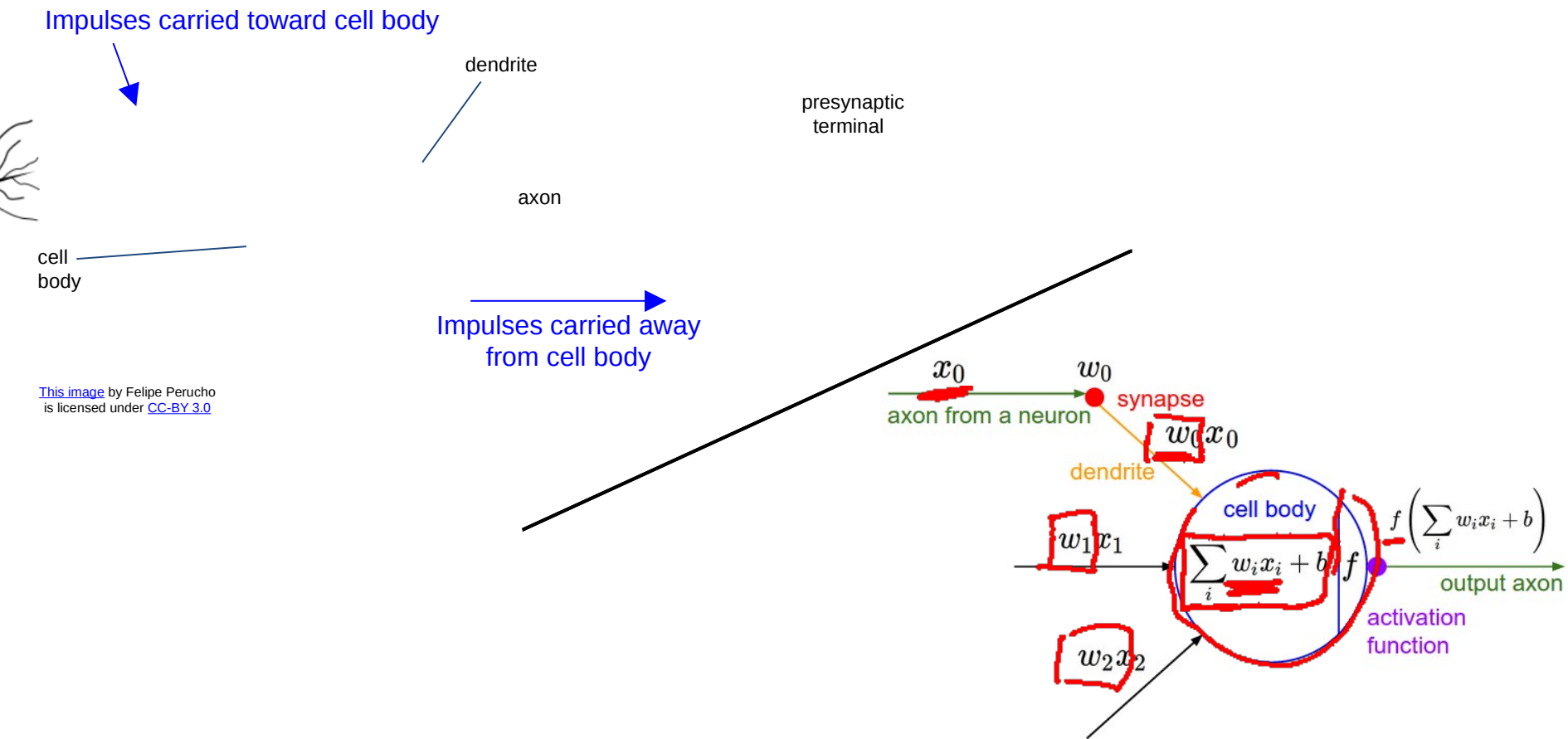


This image by [Fotis Bobolas](#) is licensed under [CC-BY 2.0](#)





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

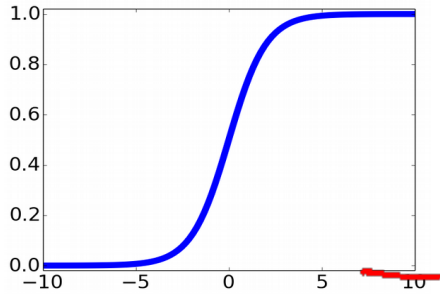


Impulses carried toward cell body



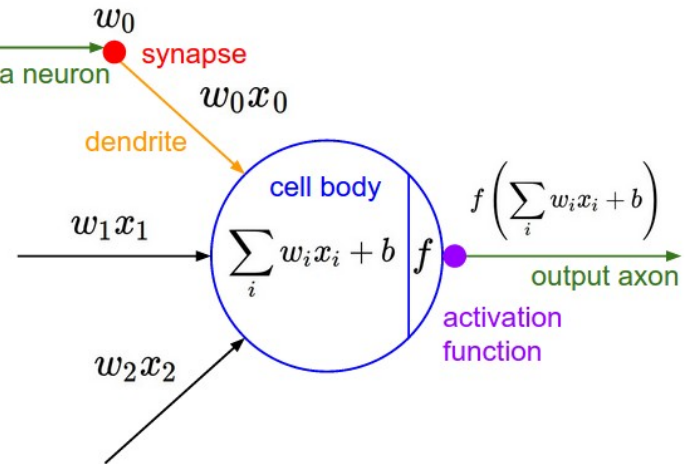
Impulses carried away from cell body

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



sigmoid activation function

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



# Be very careful with your brain analogies!

## Biological Neurons:

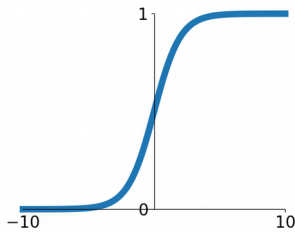
- ! ● Many different types
- Dendrites can perform complex non-linear computations
- Synapses are not a single weight but a complex non-linear dynamical system
- Rate code may not be adequate

[Dendritic Computation. London and Hausser]

# Activation functions

**Sigmoid**

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

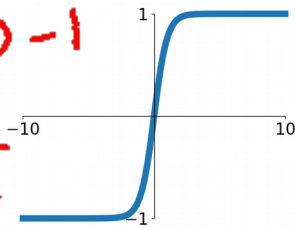


**tanh**

tanh(x)

$$2\sigma(2x) - 1$$

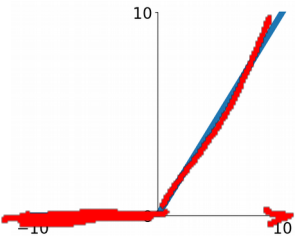
$$\frac{e^x - e^{-x}}{e^x + e^{-x}}$$



**ReLU**

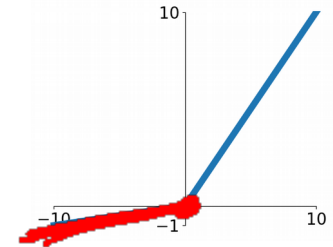
$$\max(0, x)$$

$\vec{h}$



**Leaky ReLU**

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

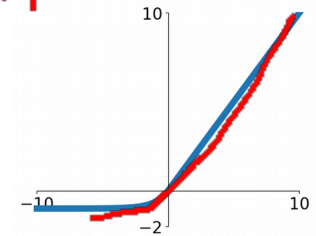


**Maxout**

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

**ELU**

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

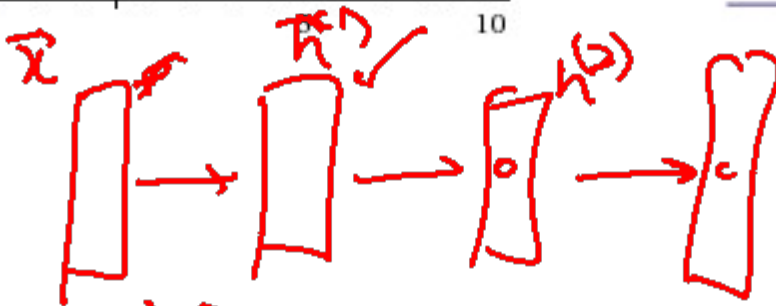
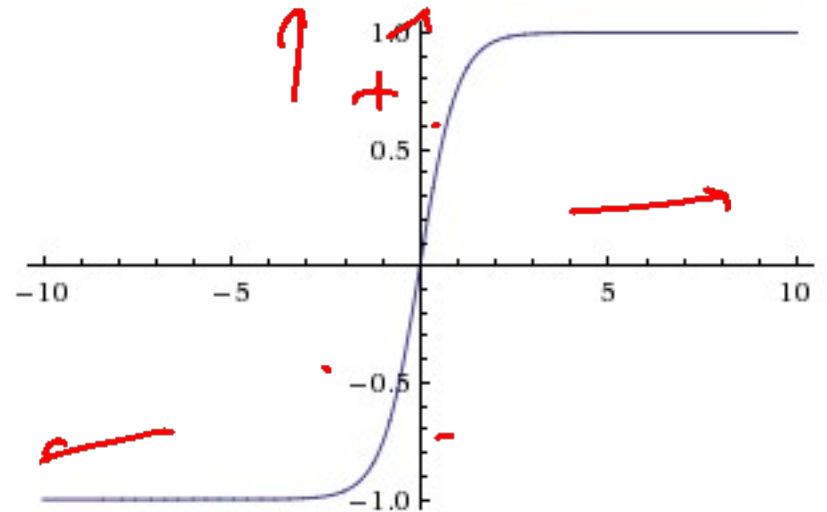
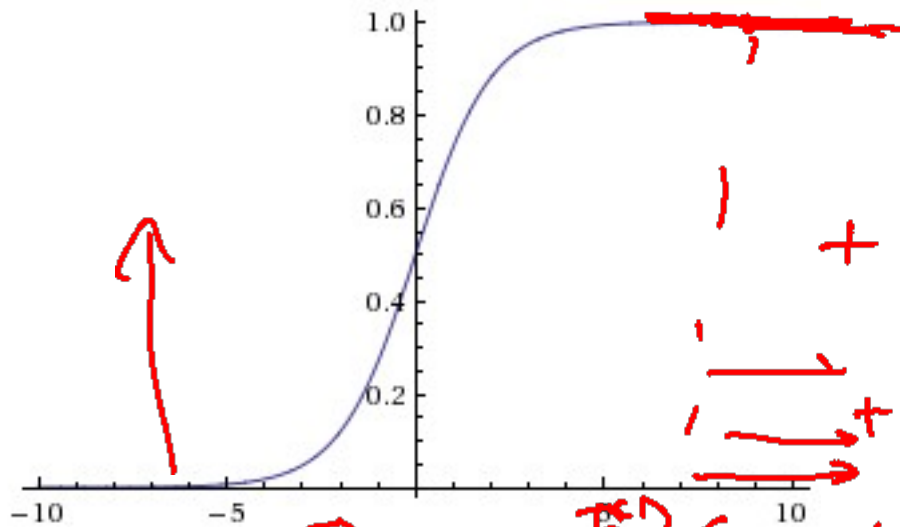


S

# Activation Functions

- sigmoid vs tanh

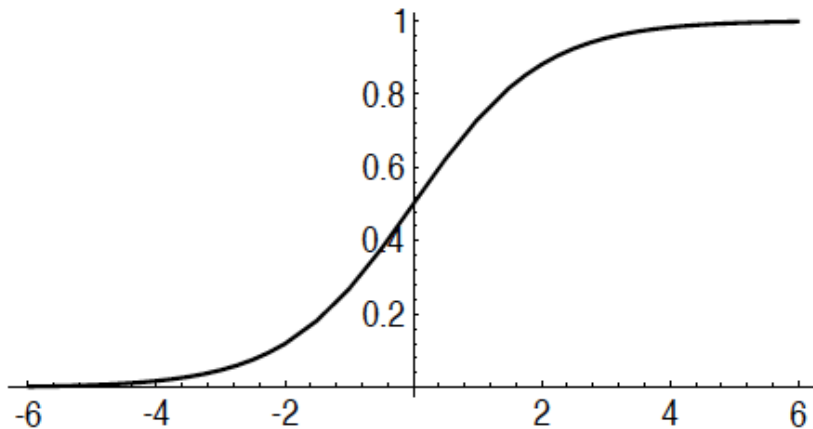
$$2\sigma(2x) - 1$$



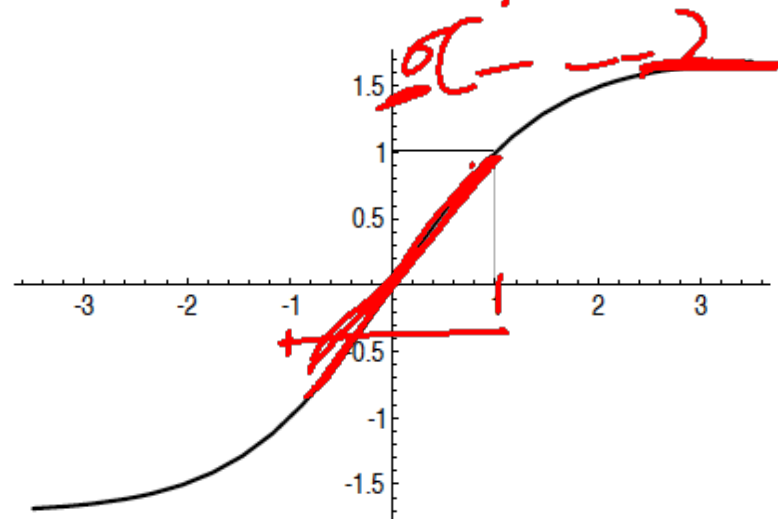
$$\sigma(\sum_i w_i h_i)$$

# A quick note

$$\sigma\left(\sum w_j h_j + b\right)$$



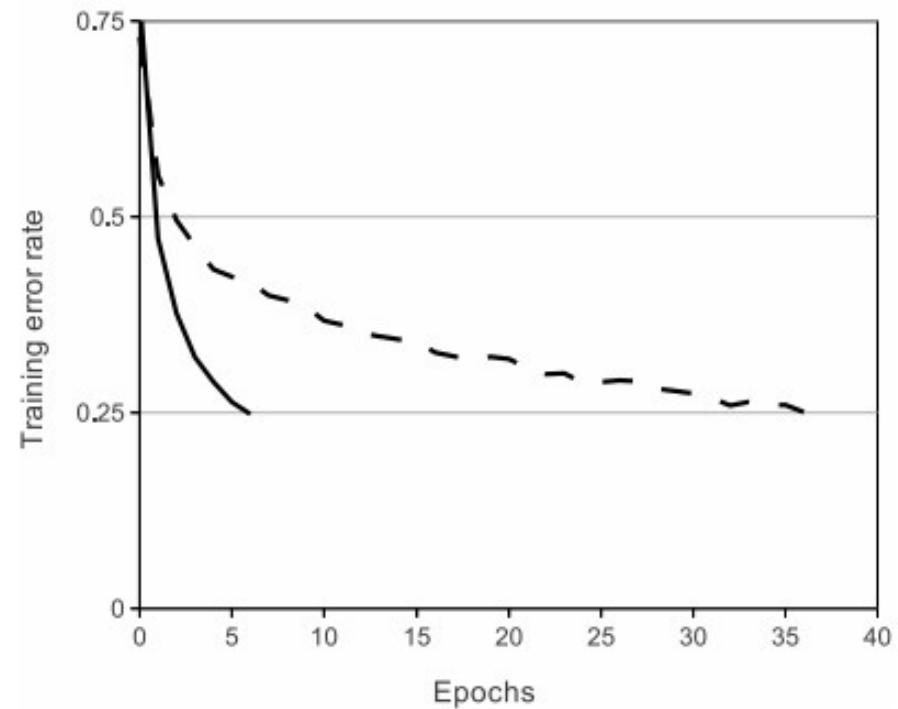
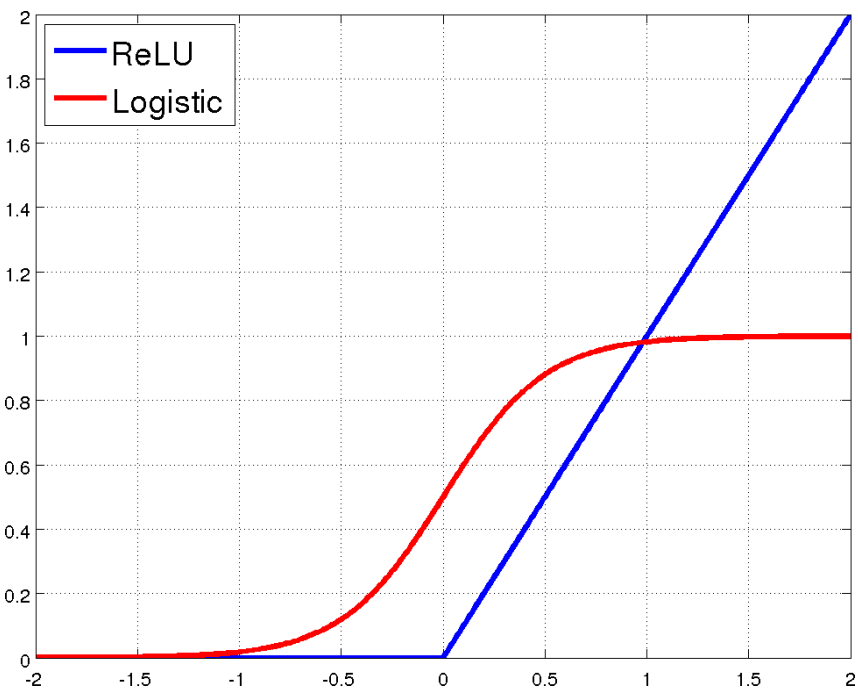
(a)



(b)

**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]



# Demo Time

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