

NEURAL NETWORK FOR MACHINE LEARNING

한재근 과장 | jahan@nvidia.com

유현곤 부장 | hryu@nvidia.com / 양한별 과장 | hanbyuly@nvidia.com



AGENDA

Background of Neural Network

Scoring / Loss / SVM / Softmax / Regularization

Neural Network Architecture

Optimization / Training

Artificial Neural Network

Activation Function

Deep Learning

Learning Representation/Features

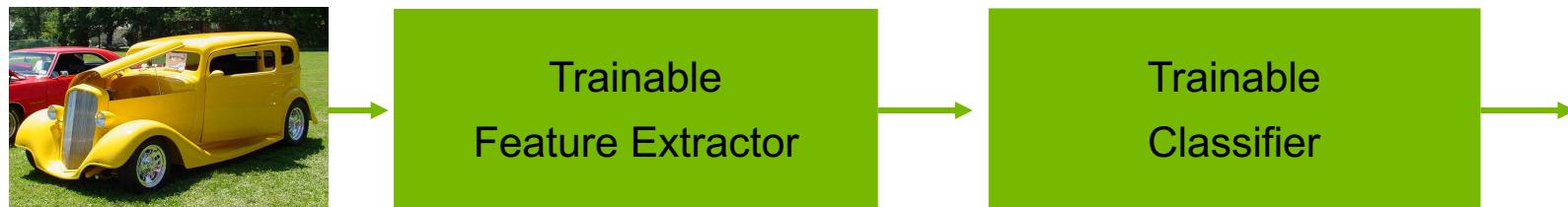
The traditional model of pattern recognition (since the late 50's)

Fixed/engineered features (or fixed kernel) + trainable classifier



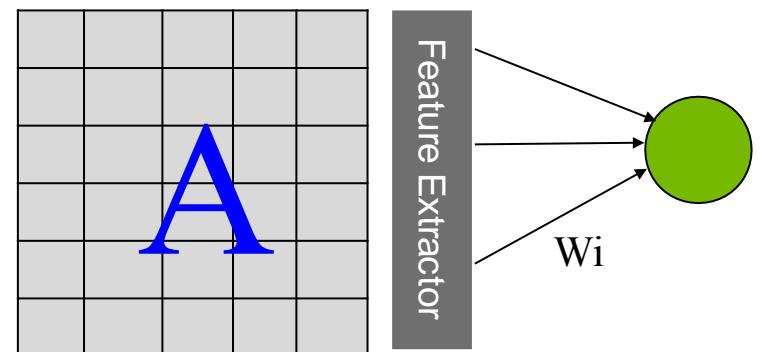
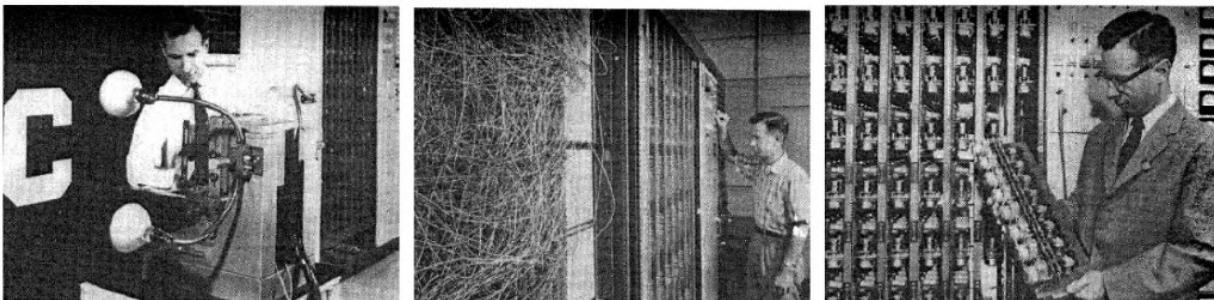
End-to-end learning / Feature learning / Deep learning

Trainable features (or kernel) + trainable classifier



This basic model has not evolved much since the 50's

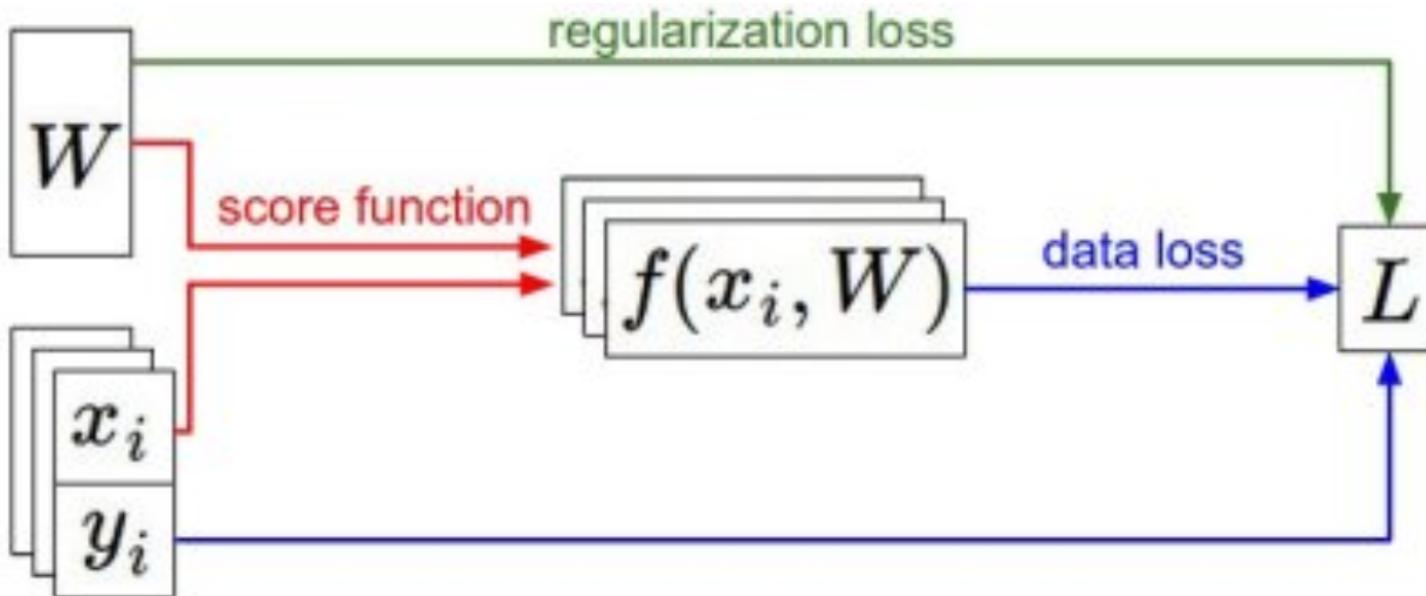
- The first learning machine: the Perceptron
 - Built at Cornell in 1960
- The Perceptron was a linear classifier on top of a simple feature extractor
- The vast majority of practical applications of ML today use glorified linear classifiers or glorified template matching.
- Designing a feature extractor requires considerable efforts by experts.

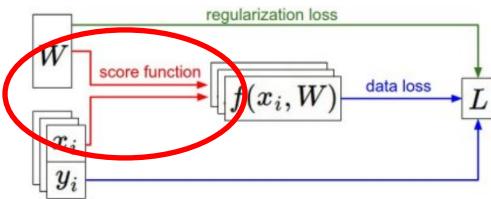


$$y = \text{sign} \left(\sum_{i=1}^N W_i F_i (X) + b \right)$$

Scoring

Summary





Inference

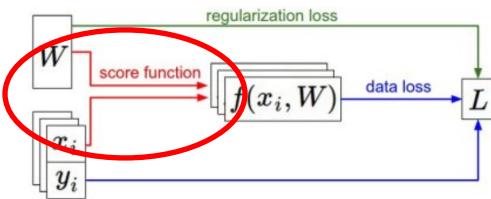
Scoring

Input Data

$$f(x_i, W, b) = Wx_i + b$$

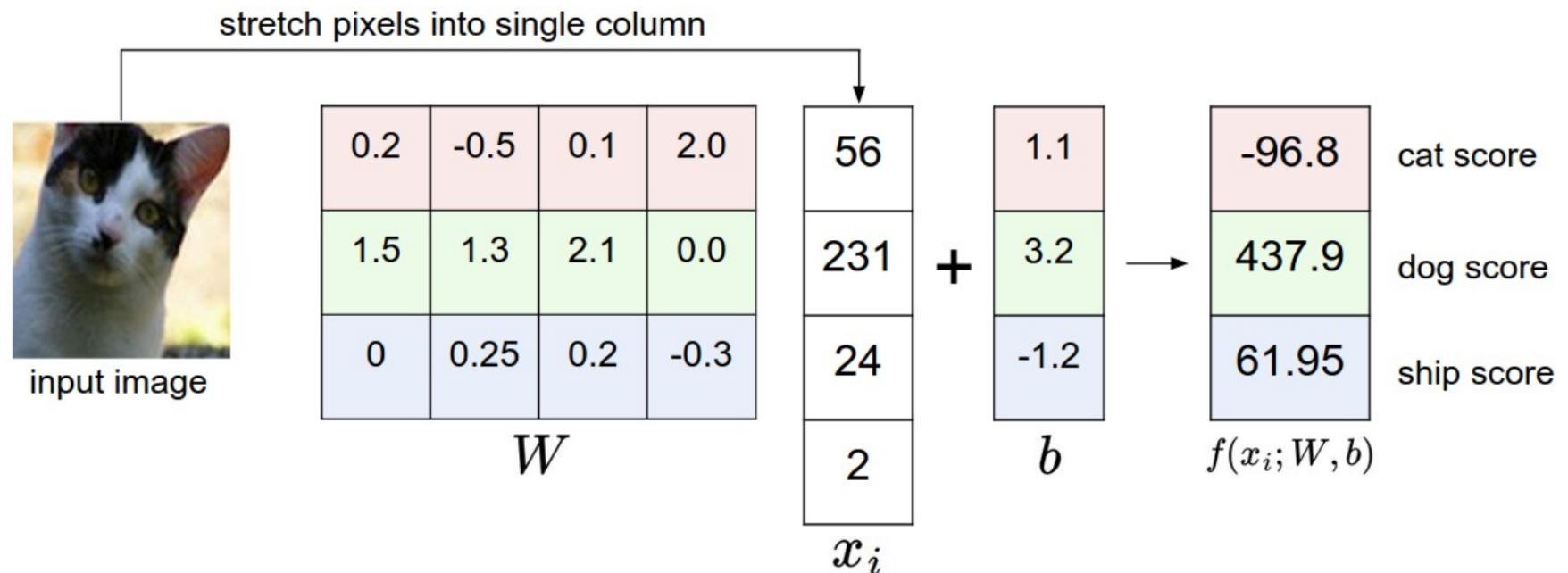
Class Score Weights Bias

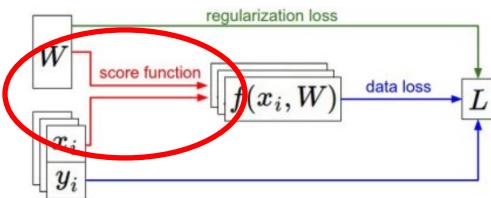
“Parameters”



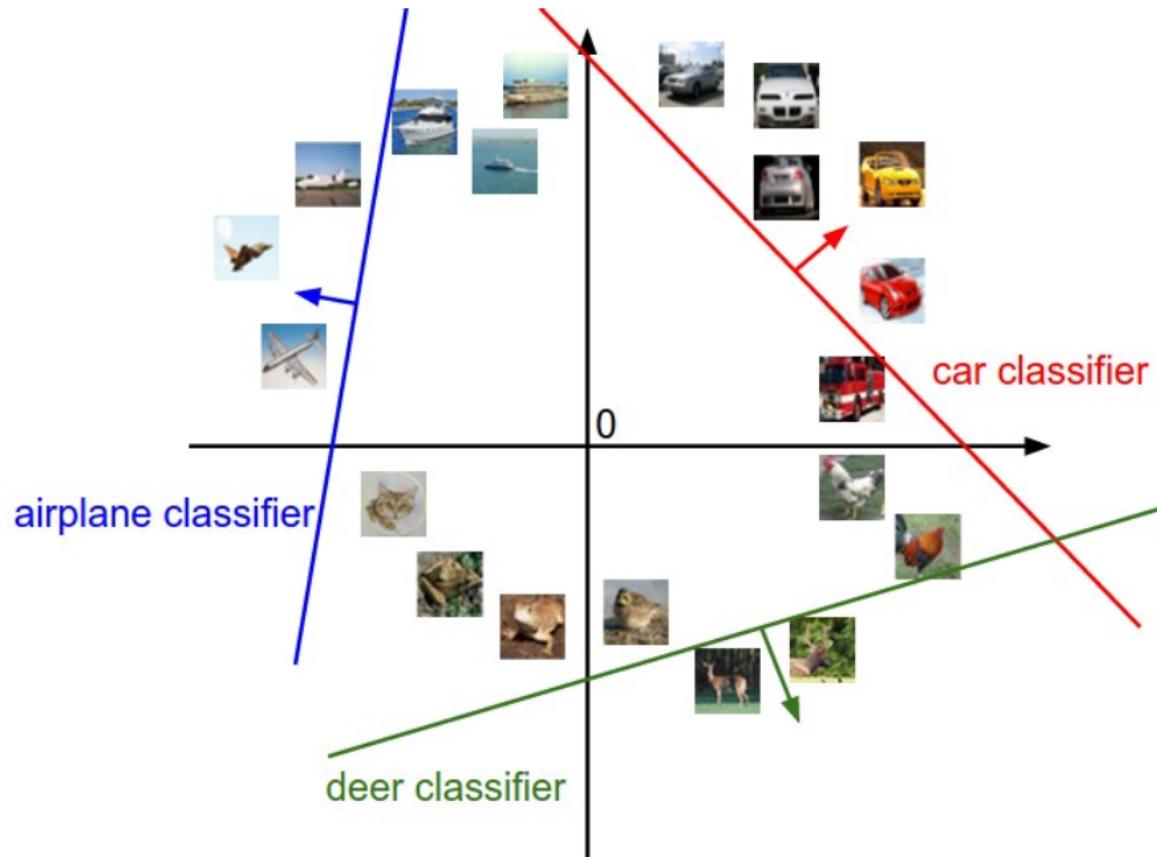
A Linear classifier for Image

$$f(x_i, W, b) = Wx_i + b$$

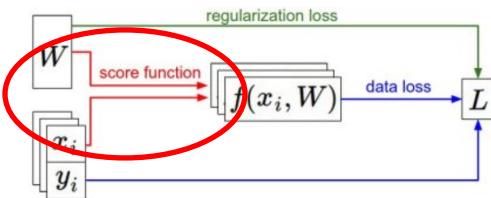




Interpretation



$$f(x_i, W, b) = Wx_i + b$$



The bias trick

$$f(x_i, W, b) = Wx_i + b \longrightarrow f(x_i, W) = Wx_i$$

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	+ 1.1
231	3.2
24	-1.2
2	

x_i

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

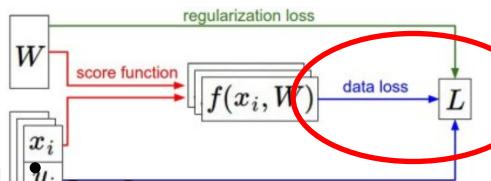
W

b

new, single W

56
231
24
2

x_i



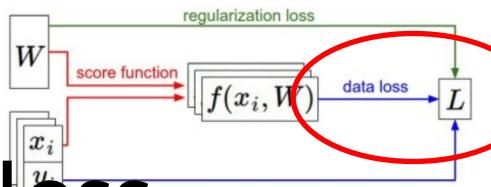
Define a loss/cost/objective function

We have **scores** and **labels**: $f(x_i, W)$ y_i

Example: $f(x_i, W) = [13, -7, 11]$

$$y_i = 0 \quad \xrightarrow{\hspace{1cm}}$$

We want to define a loss function L_i that describes how unhappy we are with these scores.

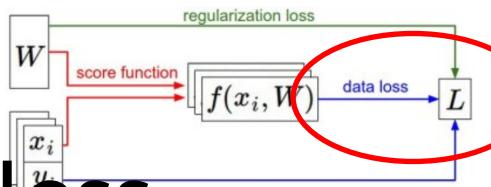


Loss example: the multiclass SVM loss

$$L_i = \sum_{j \neq y_i} \max(0, f(x_i, W)_j - f(x_i, W)_{y_i} + \Delta)$$

(One possible generalization of Binary Support Vector Machine to multiple classes)

$$L_i = C \max(0, 1 - y_i w^T x_i) + R(W)$$

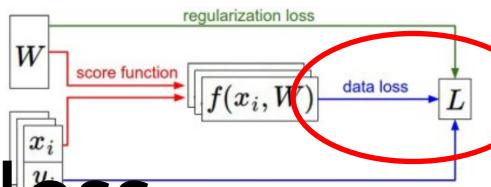


Loss example: the multiclass SVM loss

$$L_i = \sum_{j \neq y_i} \max(0, f(x_i, W)_j - f(x_i, W)_{y_i} + \Delta)$$

loss due to example i
 sum over all incorrect labels
 difference between the correct class score and incorrect class score



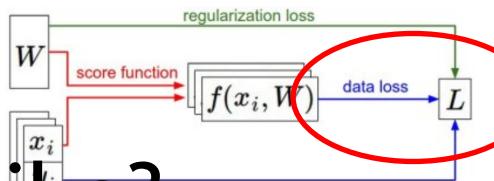


Loss example: the multiclass SVM loss

Vectorized implementation

$$L_i = \sum_{j \neq y_i} \max(0, f(x_i, W)_j - f(x_i, W)_{y_i} + \Delta)$$

```
def L_i_vectorized(x, y, W):
    scores = W.dot(x)
    margins = np.maximum(0, scores - scores[y] + 1)
    margins[y] = 0
    loss_i = np.sum(margins)
    return loss_i
```



What does the Loss function look like?

$$L_i = \sum_{j \neq y_i} \left[\max(0, w_j^T x_i - w_{y_i}^T x_i + 1) \right]$$

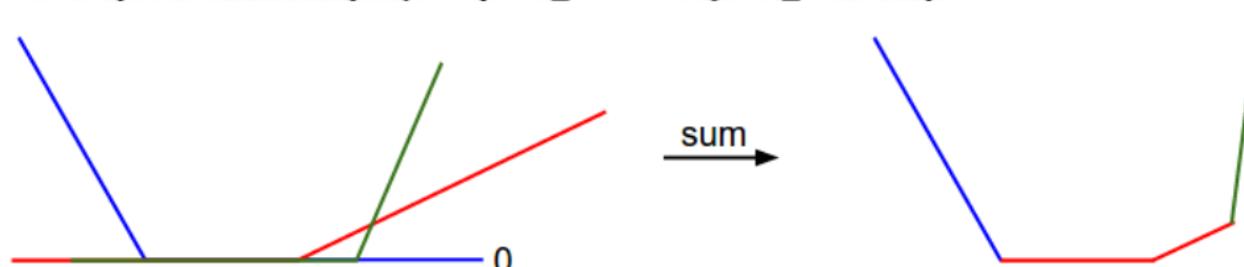
Example: 3 examples with 3 classes

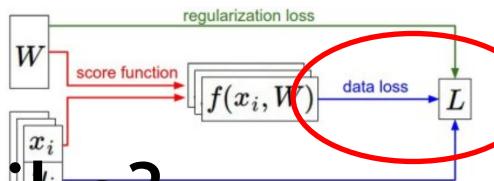
$$L_0 = \max(0, w_1^T x_0 - w_0^T x_0 + 1) + \max(0, w_2^T x_0 - w_0^T x_0 + 1)$$

$$L_1 = \max(0, w_0^T x_1 - w_1^T x_1 + 1) + \max(0, w_2^T x_1 - w_1^T x_1 + 1)$$

$$L_2 = \max(0, w_0^T x_2 - w_2^T x_2 + 1) + \max(0, w_1^T x_2 - w_2^T x_2 + 1)$$

$$L = (L_0 + L_1 + L_2)/3$$



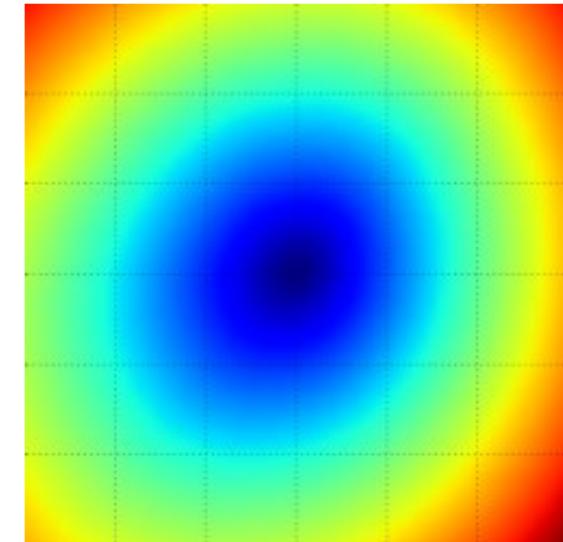
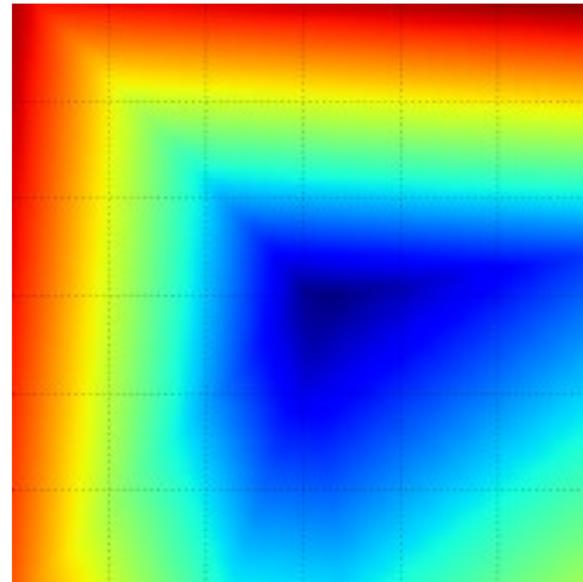


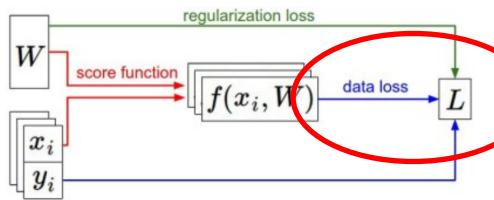
What does the Loss function look like?

$$L_i = \sum_{j \neq y_i} \max(0, f(x_i, W)_j - f(x_i, W)_{y_i} + \Delta)$$

the full data loss:

$$L_i(W + aW_1 + bW_2)$$





Softmax classifier

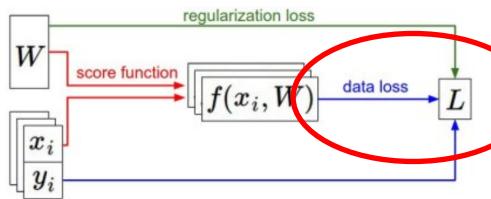
Extension of logistic regression to multiple classes

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

$P(y_i | x_i; W)$

softmax function

i.e. we're minimizing
the negative log
likelihood.



Example

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	0.0
22	0.2
-44	-0.3
56	

x_i

$+ b$

2	
---	--

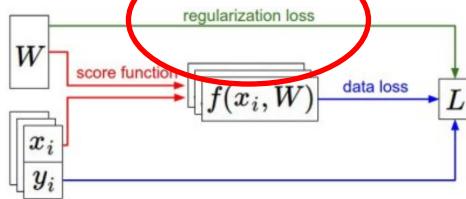
y_i

hinge loss (SVM)

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

cross-entropy loss (Softmax)

$$\begin{aligned}
 & \text{exp} \rightarrow \begin{array}{c} -2.85 \\ 0.86 \\ 0.28 \end{array} \\
 & \xrightarrow{\text{normalize (to sum to one)}} \begin{array}{c} 0.058 \\ 2.36 \\ 1.32 \end{array} \\
 & -\log(0.353) = 0.452
 \end{aligned}$$



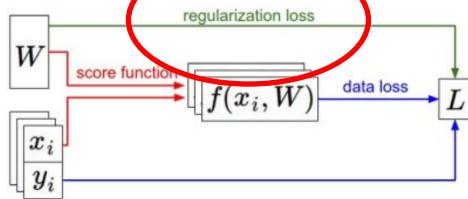
Regularization

Penalizing extreme weights

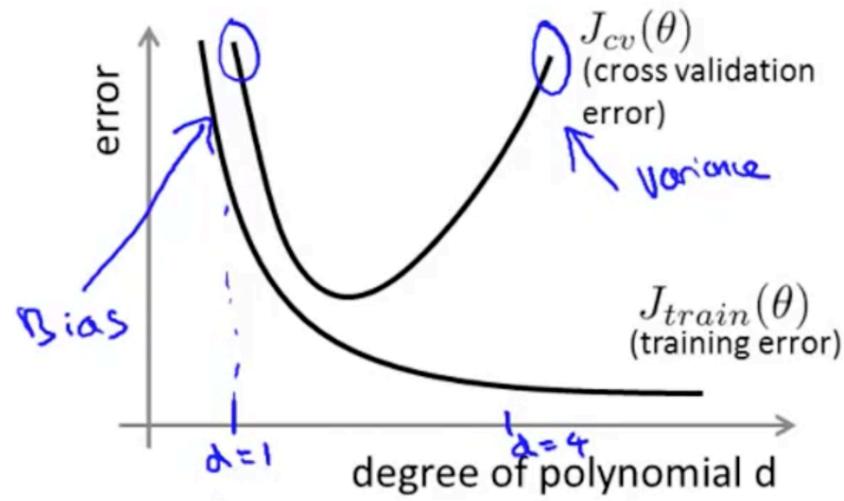
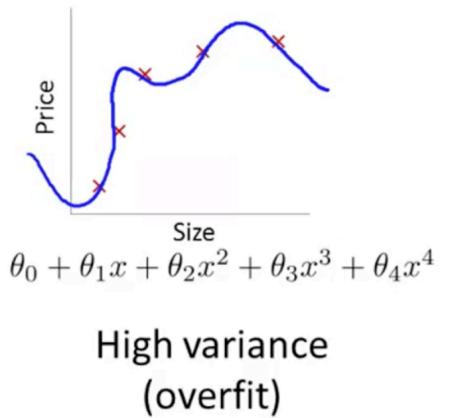
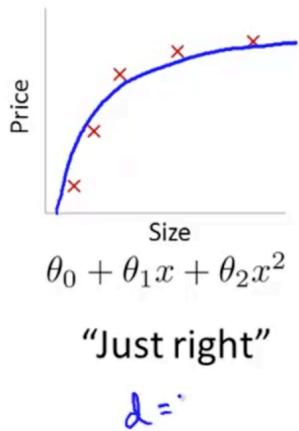
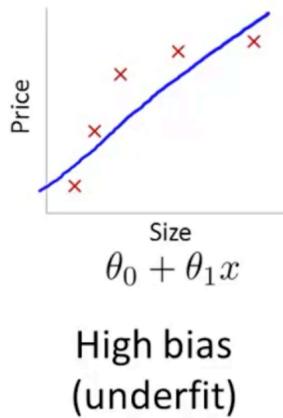
$$R(W) = \sum_k \sum_l W_{k,l}^2$$

$$L = \frac{1}{N} \sum_i \sum_{j \neq y_i} \left[\max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + \Delta) \right] + \boxed{\lambda R(W)}$$

Regularization strength



Regularization Parameter



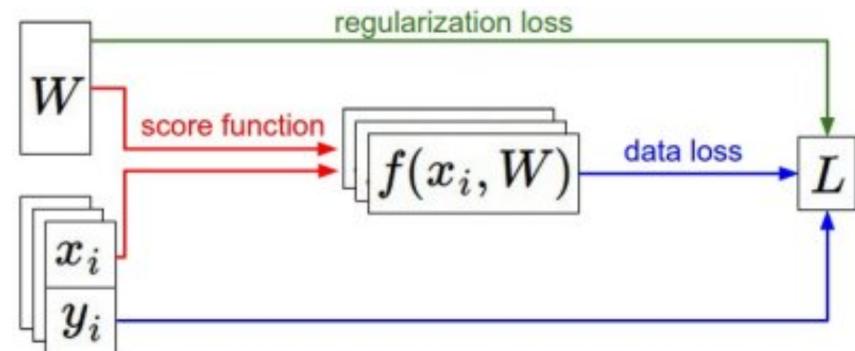
Summary

1. Score function

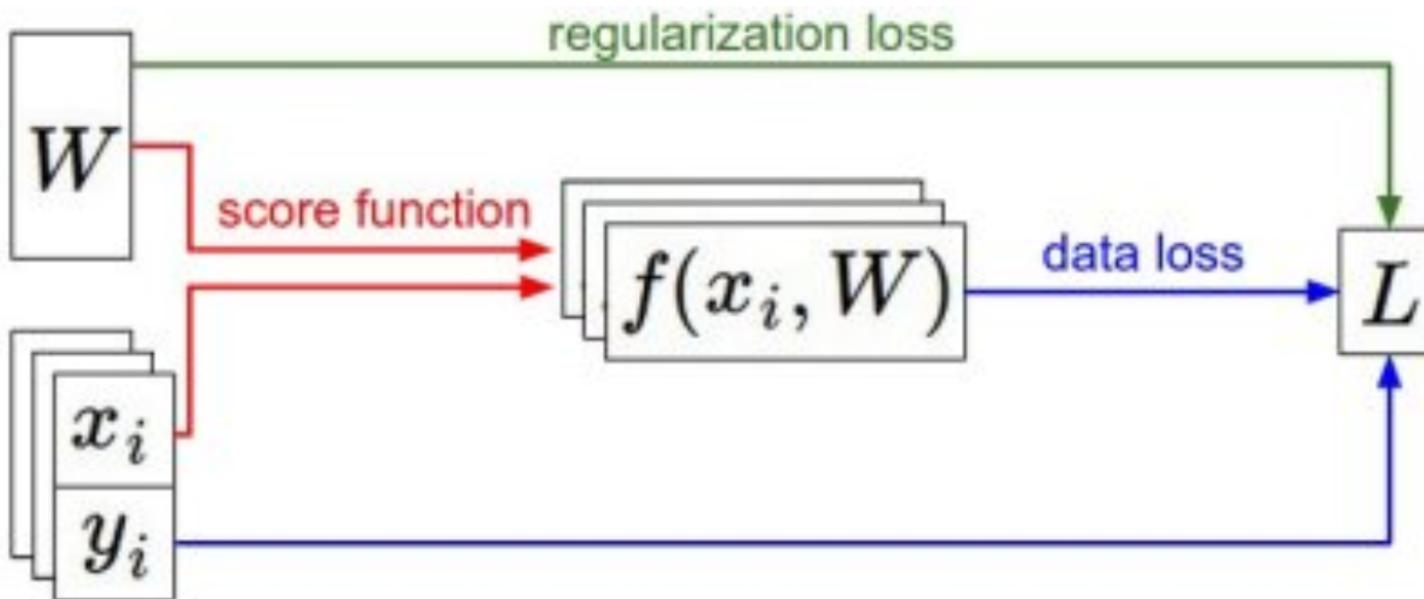
$$f(x_i, W, b) = Wx_i + b$$

2. Loss function

$$L = \frac{1}{N} \sum_i \sum_{j \neq y_i} \left[\max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + \Delta) \right] + \lambda R(W)$$



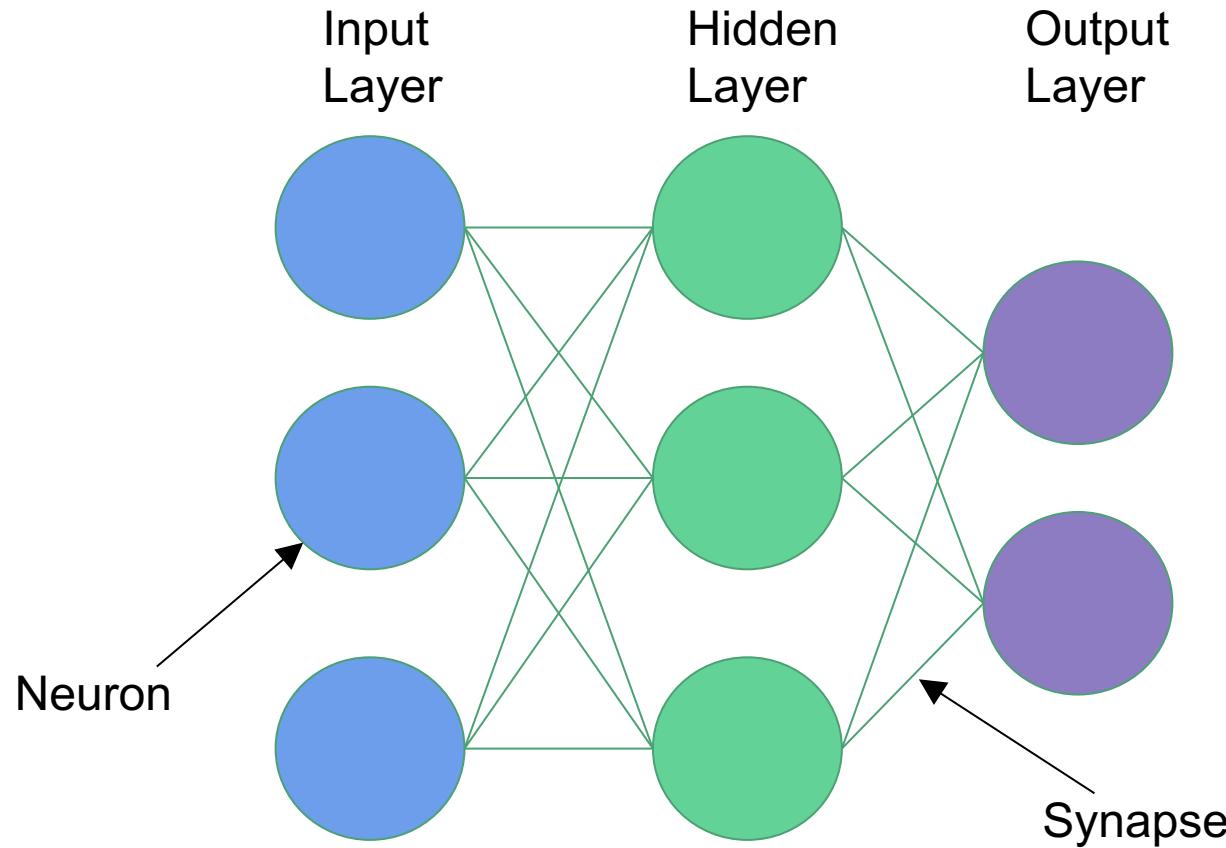
Summary



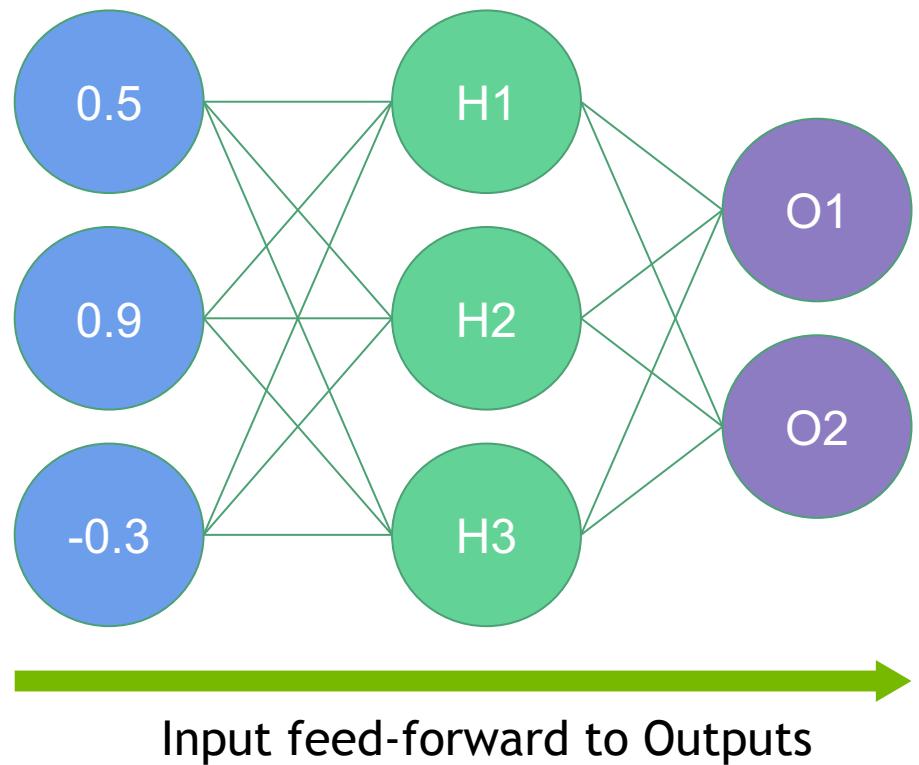
Neural Network Architecture and Optimization

Neural Network Architecture

Artificial Neural Network



Inference



H1 Weights = (1.0, -2.0, 2.0)

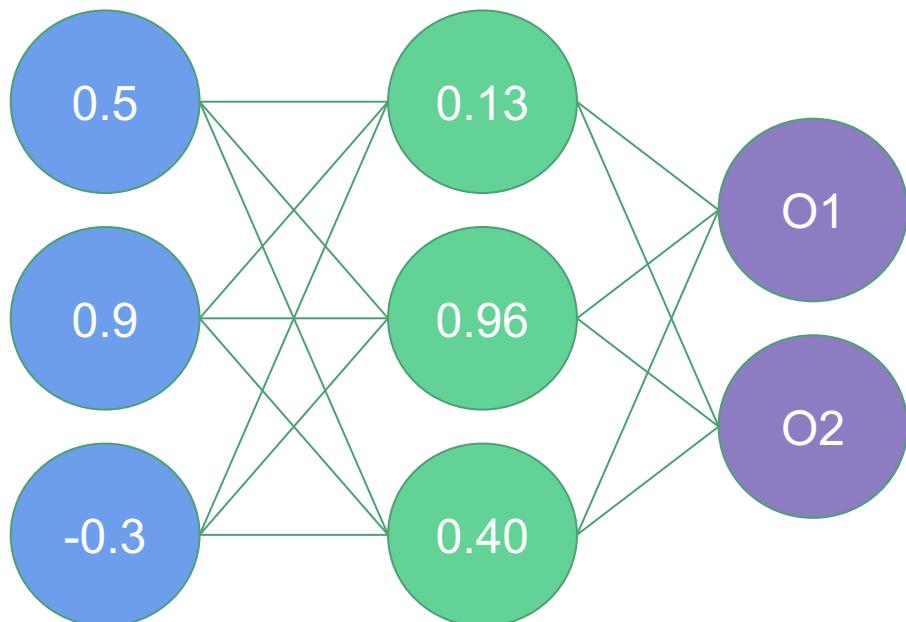
H2 Weights = (2.0, 1.0, -4.0)

H3 Weights = (1.0, -1.0, 0.0)

O1 Weights = (-3.0, 1.0, -3.0)

O2 Weights = (0.0, 1.0, 2.0)

Inference



$$\begin{aligned} H1 &= S(0.5 * 1.0 + 0.9 * -2.0 + -0.3 * 2.0) = S(-1.9) = .13 \\ H2 &= S(0.5 * 2.0 + 0.9 * 1.0 + -0.3 * -4.0) = S(3.1) = .96 \\ H3 &= S(0.5 * 1.0 + 0.9 * -1.0 + -0.3 * 0.0) = S(-0.4) = .40 \end{aligned}$$

H1 Weights = (1.0, -2.0, 2.0)

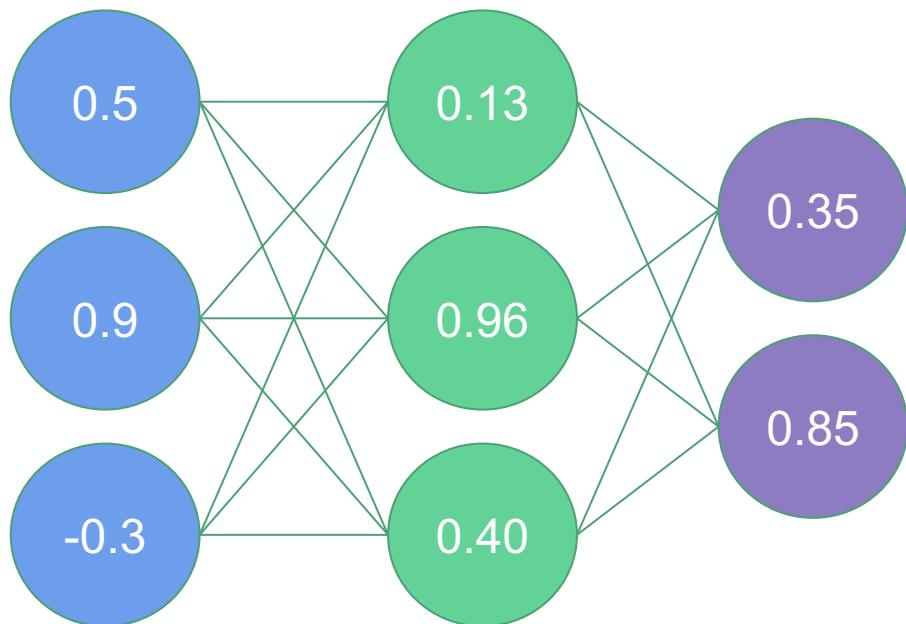
H2 Weights = (2.0, 1.0, -4.0)

H3 Weights = (1.0, -1.0, 0.0)

O1 Weights = (-3.0, 1.0, -3.0)

O2 Weights = (0.0, 1.0, 2.0)

Inference



H1 Weights = (1.0, -2.0, 2.0)

H2 Weights = (2.0, 1.0, -4.0)

H3 Weights = (1.0, -1.0, 0.0)

O1 Weights = (-3.0, 1.0, -3.0)

O2 Weights = (0.0, 1.0, 2.0)

$$O_1 = S(0.13 * -3.0 + 0.96 * 1.0 + 0.40 * -3.0) = S(-0.63) = .35$$

$$O_2 = S(0.13 * 0.0 + 0.96 * 1.0 + 0.40 * 2.0) = S(1.76) = .85$$

Matrix Formulation

H1 Weights = (1.0, -2.0, 2.0)

H2 Weights = (2.0, 1.0, -4.0)

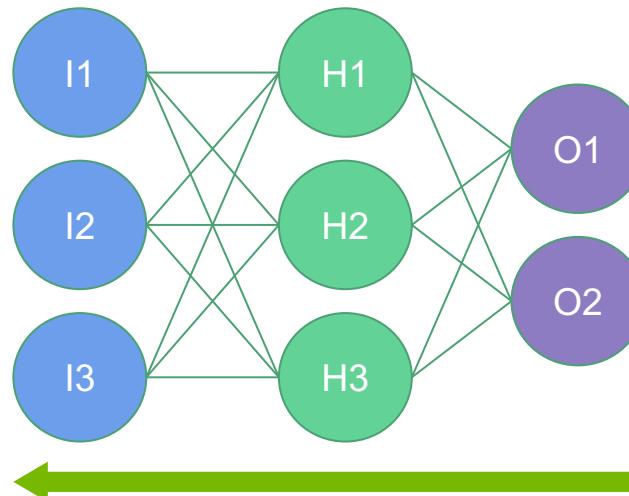
H3 Weights = (1.0, -1.0, 0.0)

Hidden Layer Weights Inputs Hidden Layer Outputs

$$S(\begin{array}{|c|c|c|} \hline 1.0 & -2.0 & 2.0 \\ \hline 2.0 & 1.0 & -4.0 \\ \hline 1.0 & -1.0 & 0.0 \\ \hline \end{array}) * \begin{array}{|c|} \hline 0.5 \\ \hline 0.9 \\ \hline -0.3 \\ \hline \end{array}) = S(\begin{array}{|c|c|c|} \hline -1.9 & 3.1 & -0.4 \\ \hline \end{array}) = \begin{array}{|c|c|c|} \hline .13 & .96 & 0.4 \\ \hline \end{array}$$

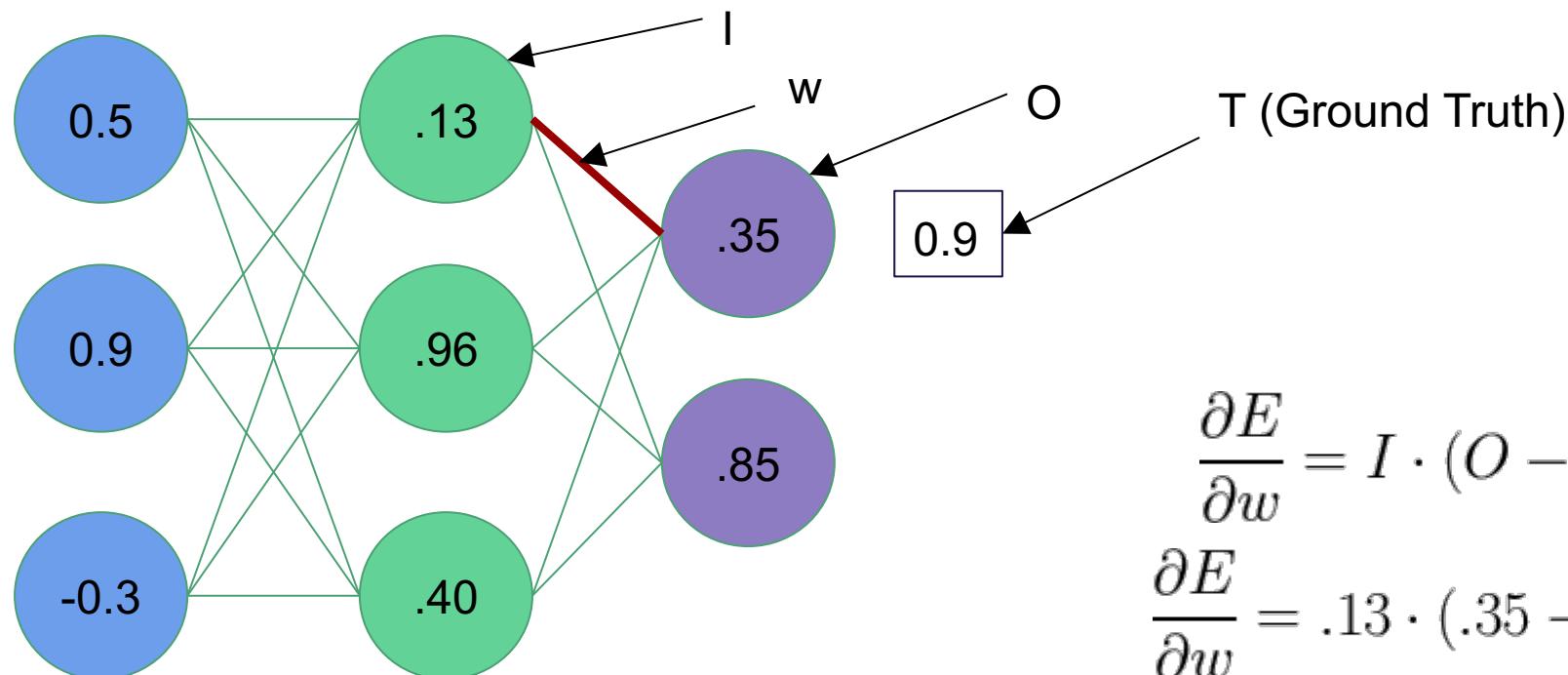
Backpropagation

- Problem: Which weights should be updated and by how much?
 - Insight: Use the derivative of the error with respect to weight to assign “blame”



backward to input to update parameters with loss

Backpropagation Example



$$\frac{\partial E}{\partial w} = I \cdot (O - T) \cdot O \cdot (1 - O)$$

$$\frac{\partial E}{\partial w} = .13 \cdot (.35 - .9) \cdot .35 \cdot (1 - .35)$$

Gradient Descent

- Gradient Descent minimizes the neural network's error
 - At each time step the error of the network is calculated on the training data
 - Then the weights are modified to reduce the error
- Gradient Descent terminates when
 - The error is sufficiently small
 - The max number of time steps has been exceeded

Gradient Descent

Getting to a minimum

The derivative of a function in 1 dimension is:

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

In n dimensions, the **gradient** is the vector of (partial) derivatives.

Gradient Descent

```
# Vanilla Gradient Descent

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

Dynamics of Gradient Descent

Intuition

SVM

pull some weights up and some down

$$L = \frac{1}{N} \sum_i \sum_{j \neq y_i} [\max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + \Delta)] + \lambda \sum_k \sum_l W_{k,l}^2$$

Softmax

$$L = \frac{1}{N} \sum_i -\log \left(\frac{e^{f_{y_i}}}{\sum_j e^{f_j}} \right) + \lambda \sum_k \sum_l W_{k,l}^2$$

always pull the weights down

Gradient intuition

$$f(x, y) = xy \quad \rightarrow \quad \frac{\partial f}{\partial x} = y \quad \frac{\partial f}{\partial y} = x$$

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

$$f(x + h) = f(x) + h \frac{df(x)}{dx}$$

Example: $x = 4, y = -3$. $\Rightarrow f(x, y) = -12$

$$\boxed{\frac{\partial f}{\partial x} = -3}$$

$$\boxed{\frac{\partial f}{\partial y} = 4}$$

partial derivatives

$$\boxed{\nabla f = \left[\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right]}$$

gradient

If I increase x by h , how would the input of f change?

Forward pass and backpropagation intuition

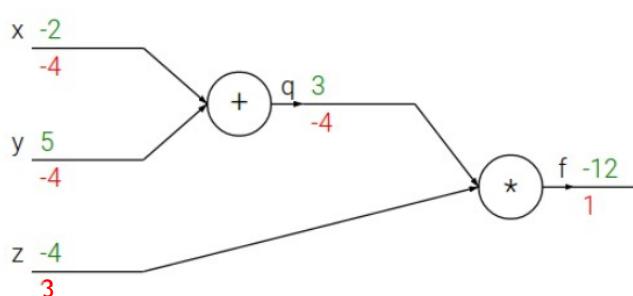
Compound expressions: $f(x, y, z) = (x + y)z$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Chain rule:

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

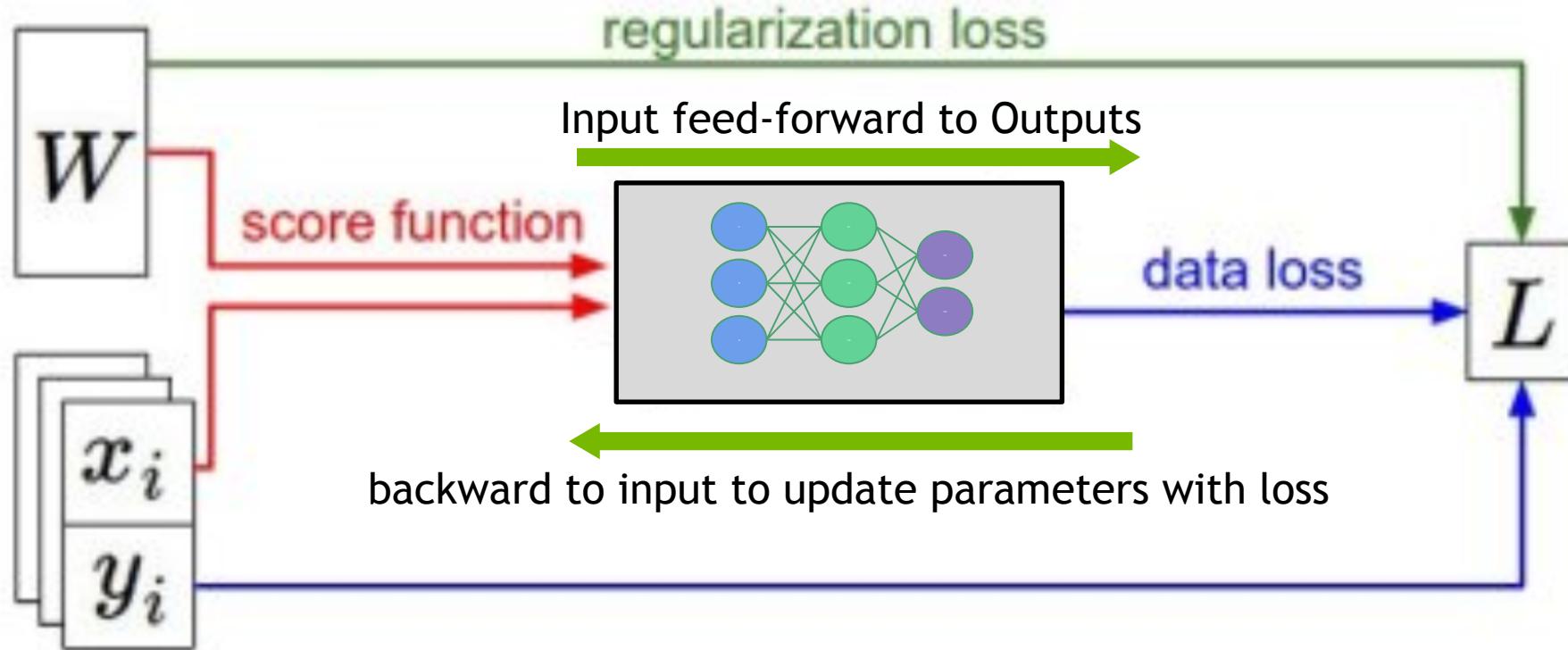


```
# set some inputs
x = -2; y = 5; z = -4

# perform the forward pass
q = x + y # q becomes 3
f = q * z # f becomes -12
```

```
# perform the backward pass (backpropagation) in reverse order:
# first backprop through f = q * z
dfdq = q # df/fz = q, so gradient on z becomes 3
dfdz = z # df/fz = z, so gradient on q becomes -4
# now backprop through q = x + y
dfdx = 1.0 * dfdq # dq/dx = 1. And the multiplication here is the chain rule!
dfdy = 1.0 * dfdq # dq/dy = 1
```

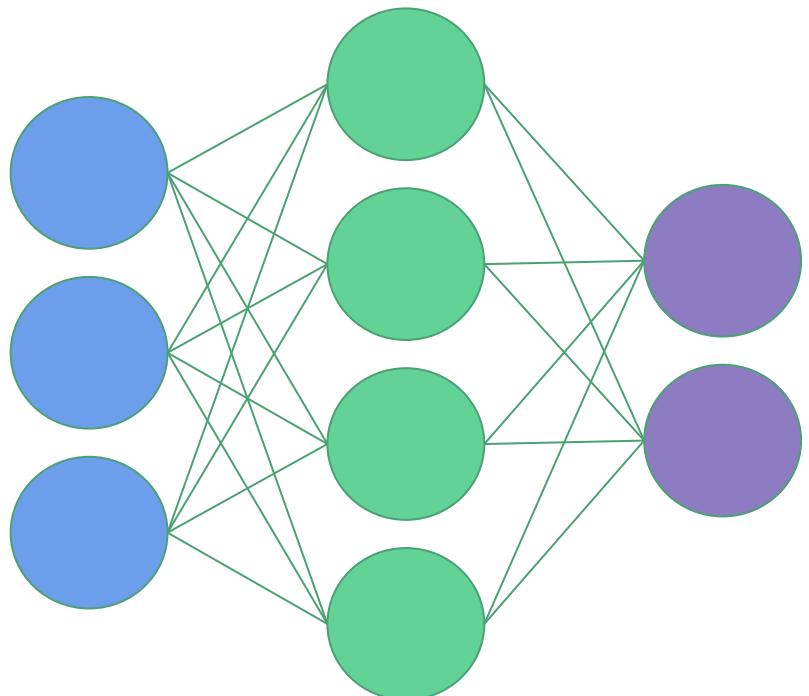
Summary



Artificial Neural Network

ANN Architecture

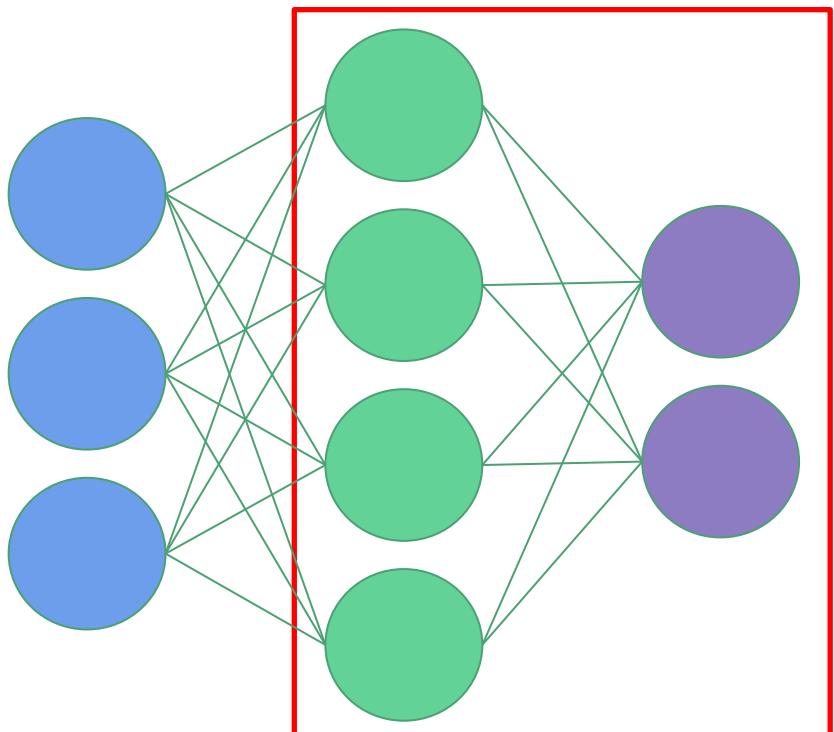
Fully Connected Layers



2-layer neural nets or
1-hidden neural nets

ANN Architecture

Trainable Parameters

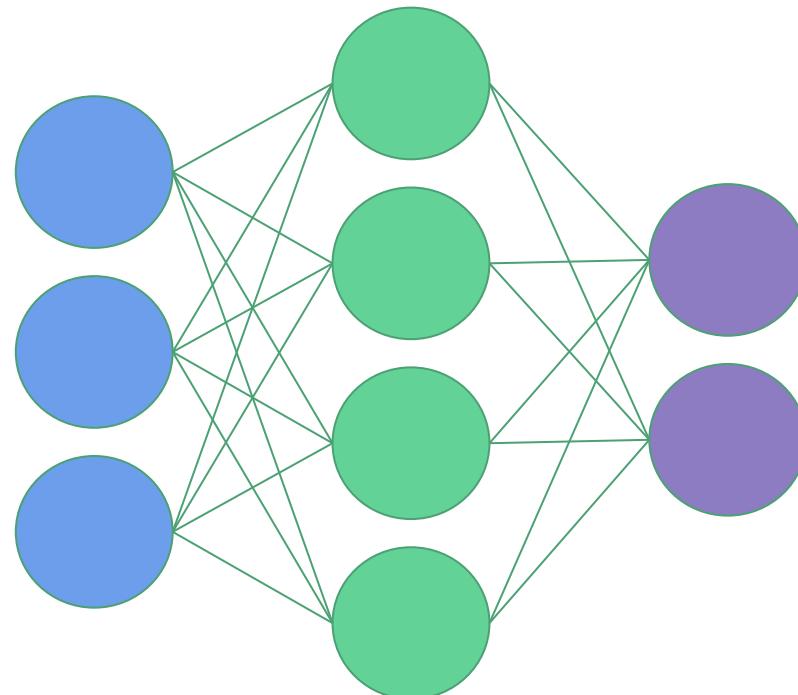


Number of Neurons: $4 + 2 = 6$
Number of Weights: $[4 \times 3 + 2 \times 4] = 20$
Number of Parameters: $6 + 20 = 26$

```
# forward-pass of a 3-layer neural network
f = lambda x: np.max(0,x)      # activation function (ReLU)
x = np.random.rand(3,1)         # random input vector of three numbers (3)
h1 = f(np.dot(W1, x) + b1)    # calculate hidden layer activations (4x1)
out = f(np.dot(W2, h1) + b2)   # output neuron (2x1)
```

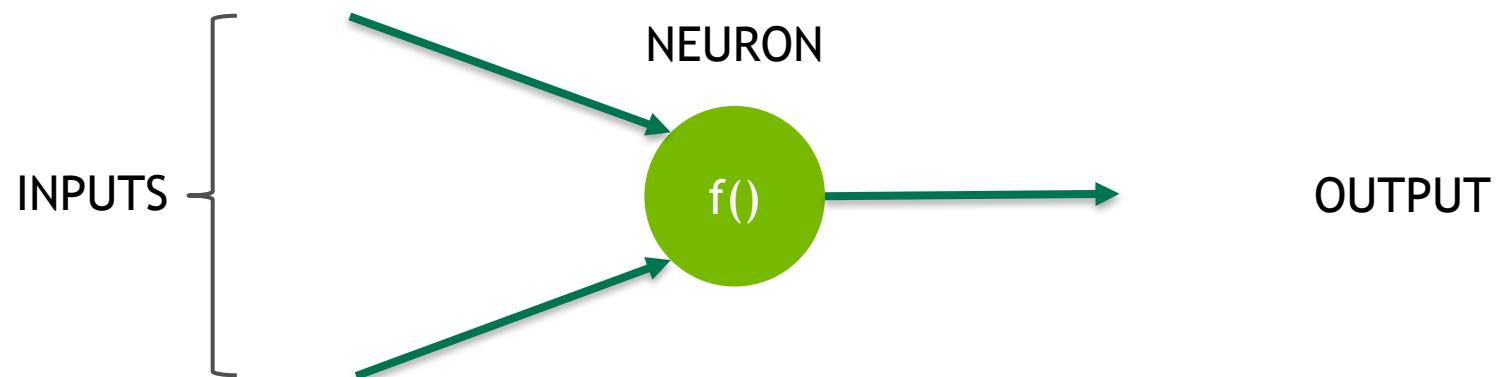
ANN Architecture

Trainable Parameters

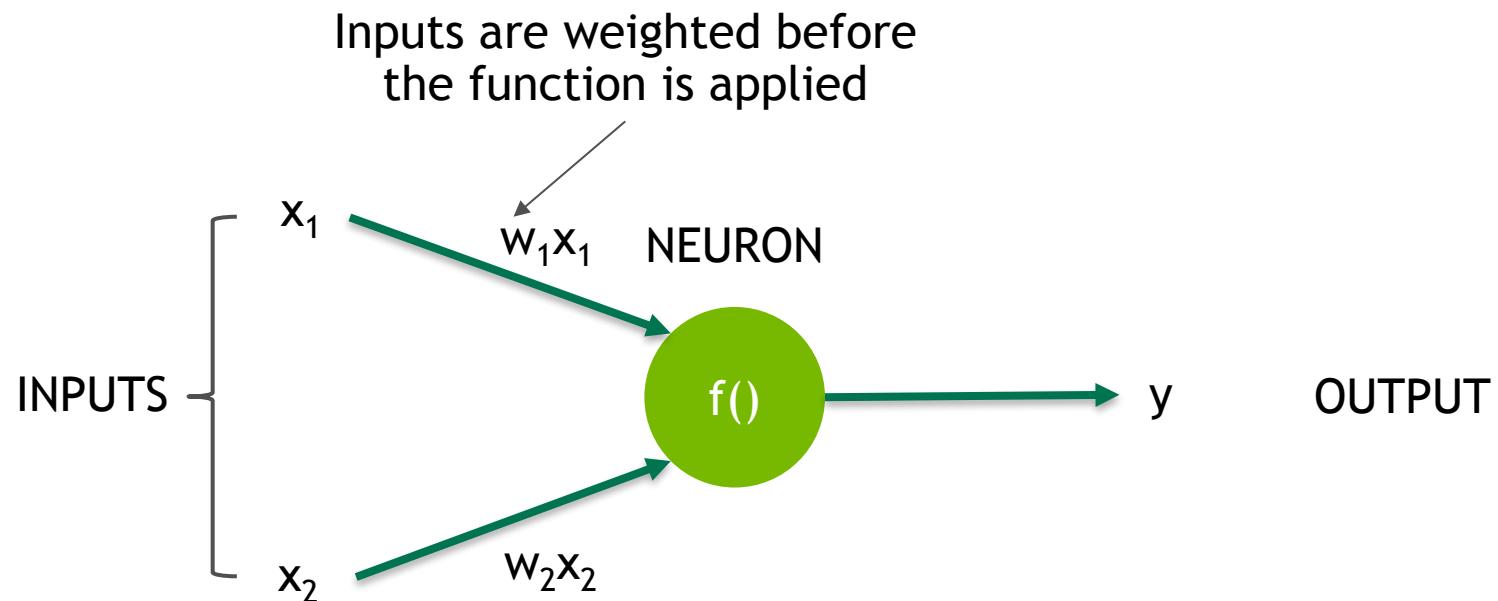


Modern CNNs: ~10 million neurons
Human visual cortex: ~5 billion neurons

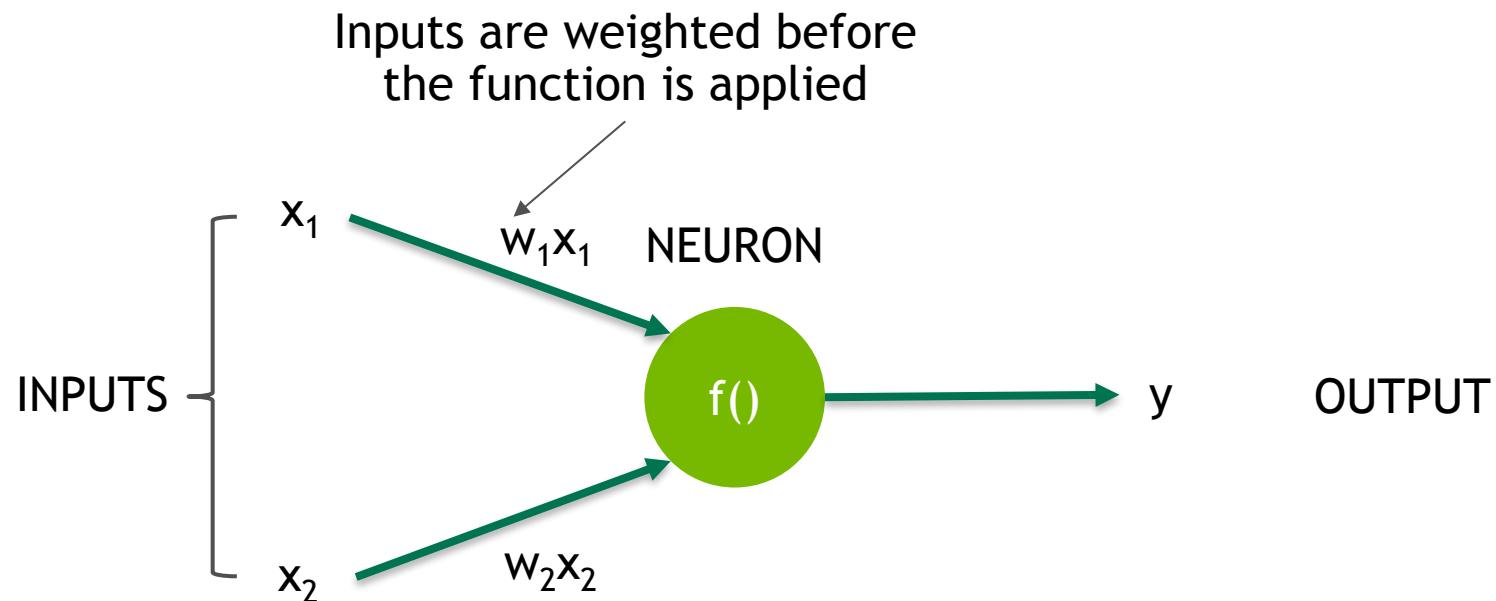
The artificial neuron



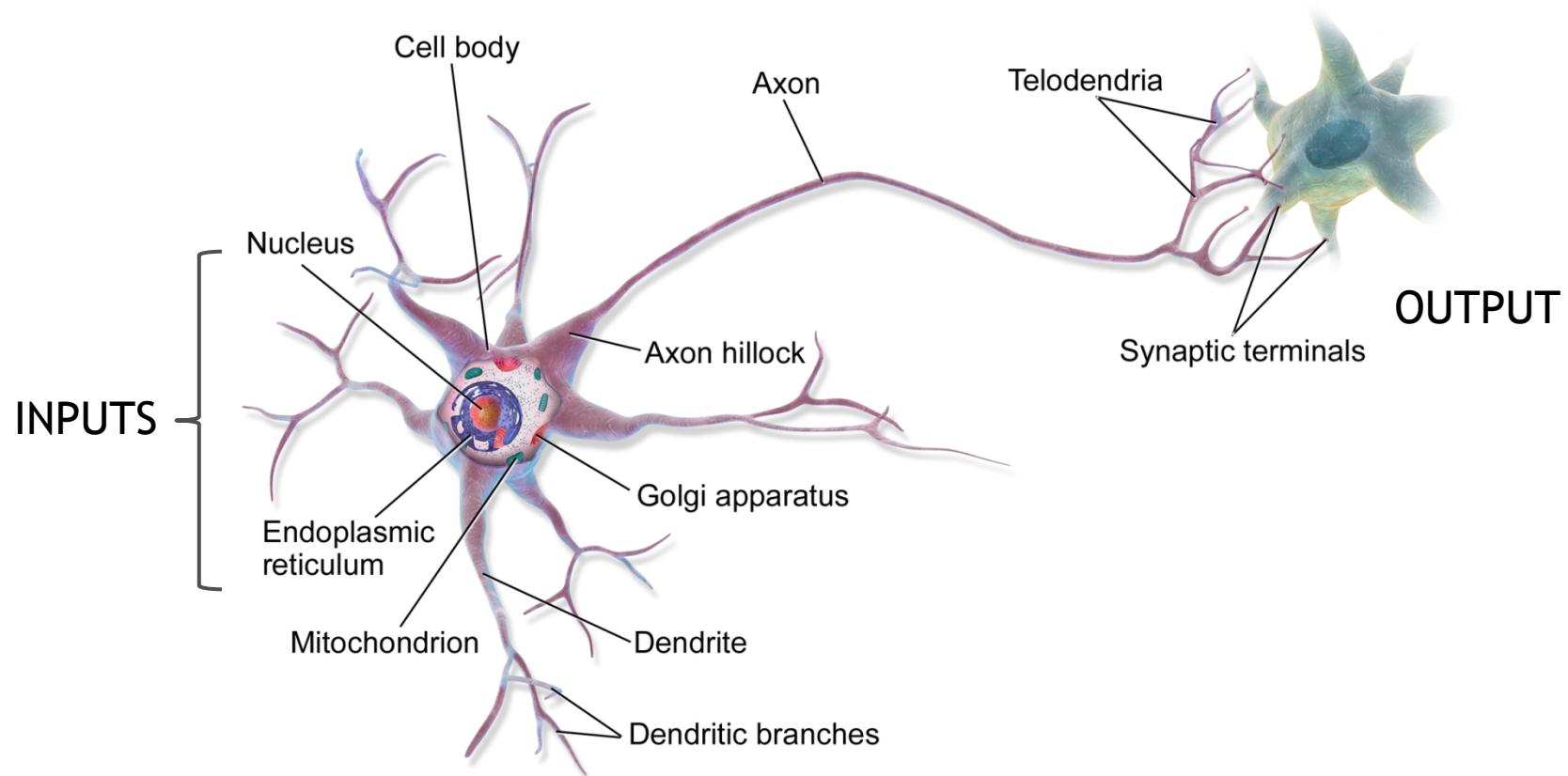
The artificial neuron



The artificial neuron



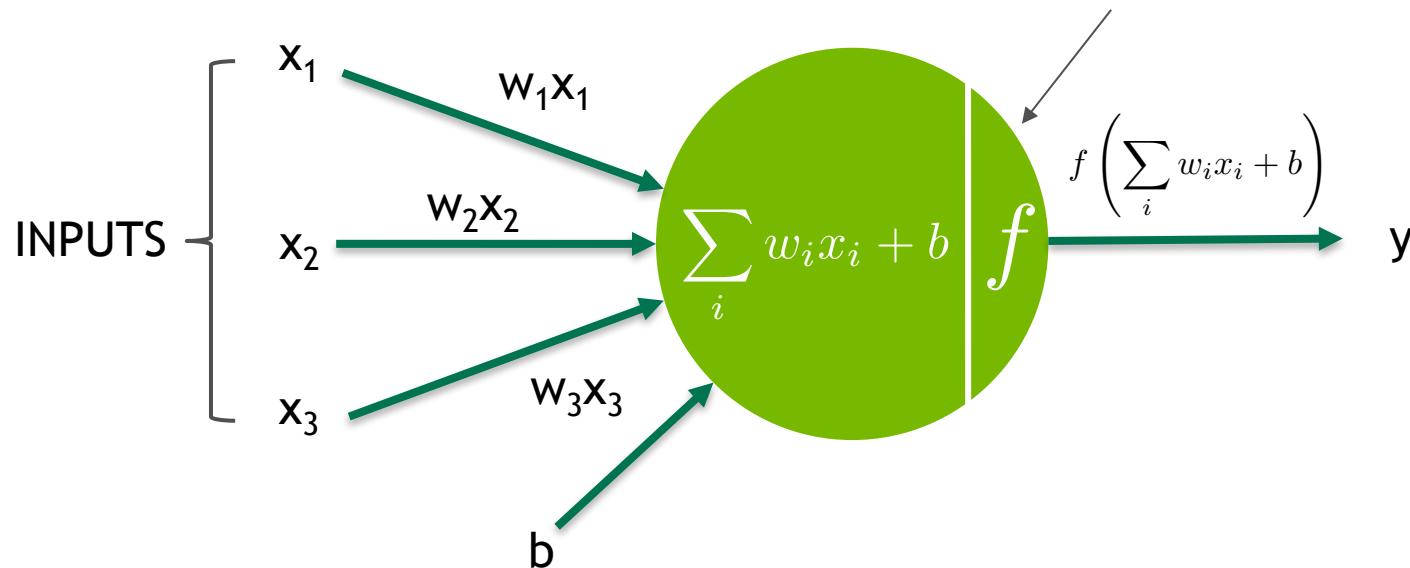
Biological Inspiration



The artificial neuron

More detailed view

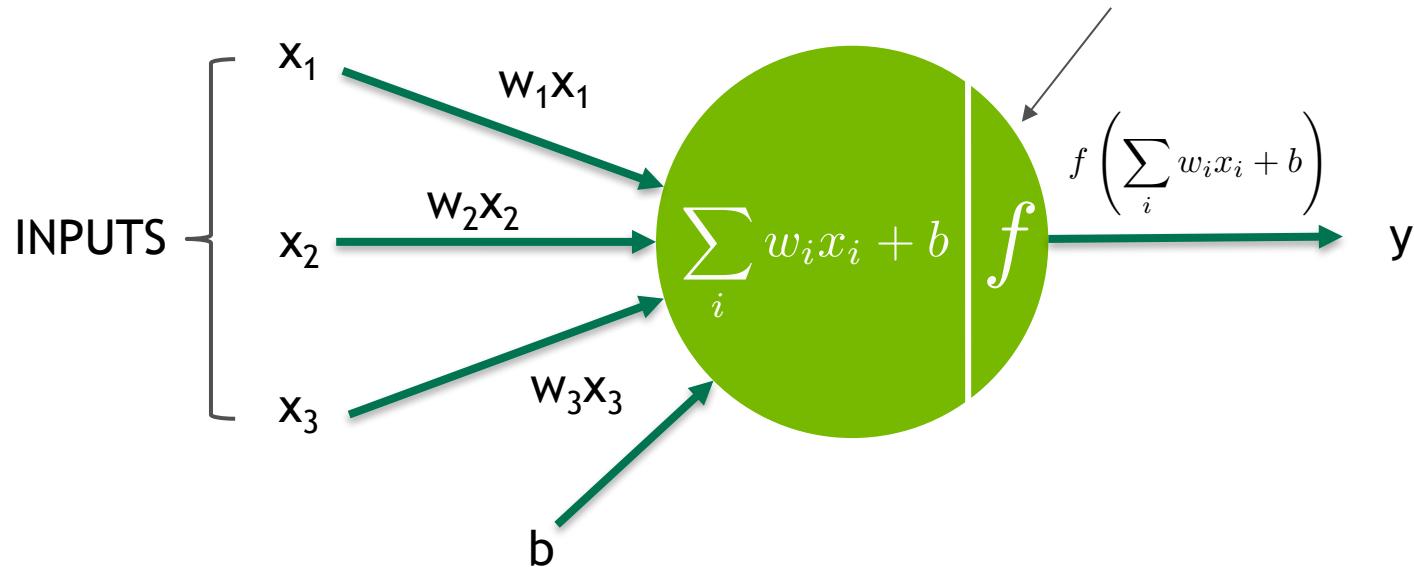
Activation function f



The artificial neuron

More detailed view

Activation function f

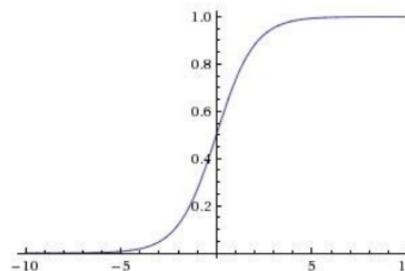


Example Python implementation:

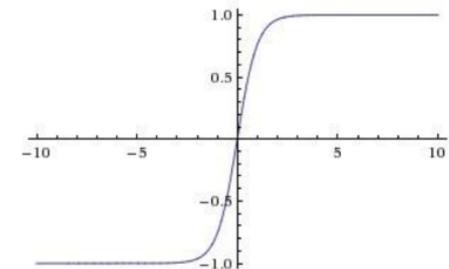
```
1 class Neuron:
2
3     def neuron_fire(self, inputs):
4         """ assume inputs and weights are 1-D arrays and bias is a number """
5         cell_body_sum = np.sum(inputs * self.weights) + self.bias
6         activation = np.max(0,cell_body_sum) # ReLU activation function
7         return activation
```

Activation functions

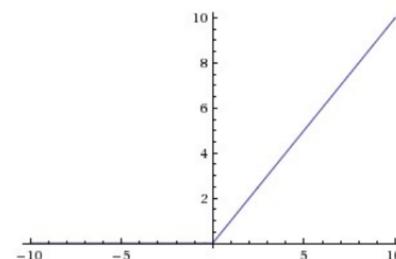
- ▶ Activation Functions are applied to the inputs at each neuron
- ▶ A range of non-linear functions have been used
- ▶ Non-linearity is important because a multi-layer ANN with linear activations would still only be able to learn composite linear functions



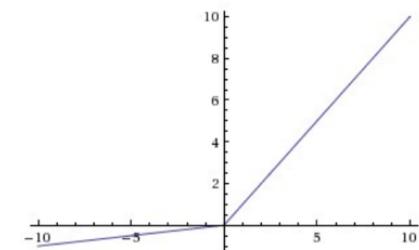
Sigmoid



$\tanh(x)$



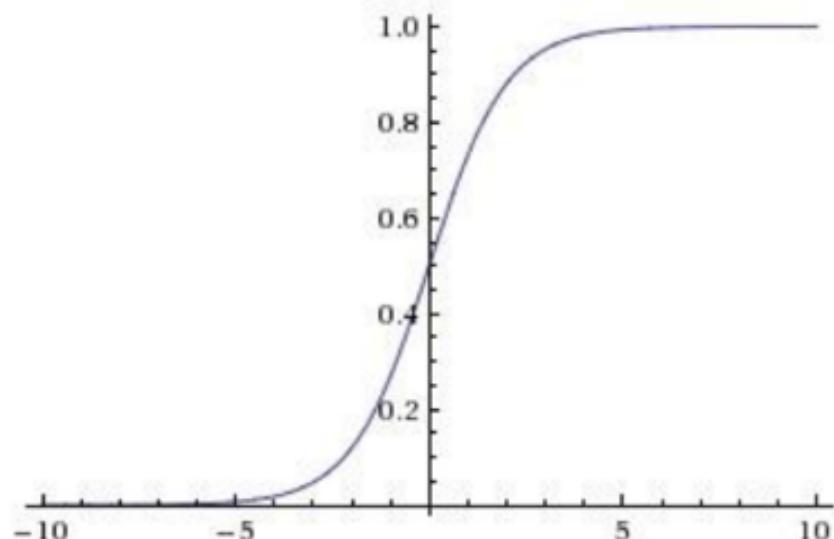
ReLU



Leaky ReLU

Activation functions

Sigmoid



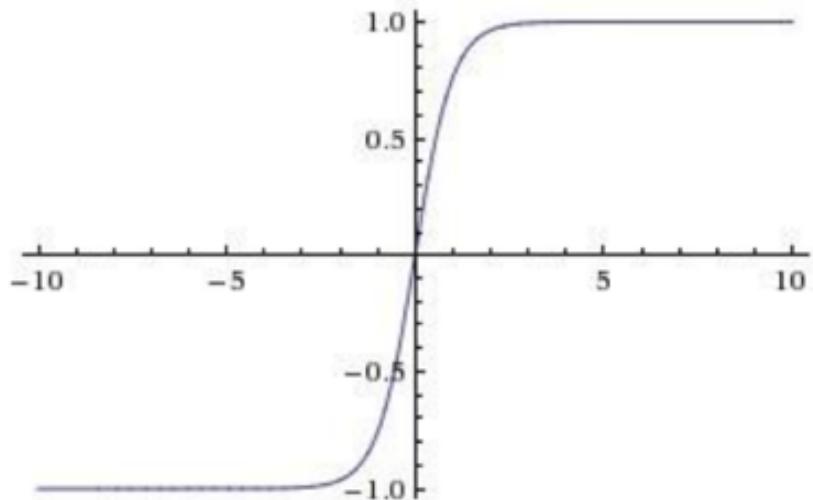
$$\sigma(x) = 1/(1 + e^{-x})$$

- ▶ Squashes numbers to the range [0,1]
- ▶ Historically popular, but they have big practical problems that prevent network training

Sigmoid

Activation functions

Hyperbolic Tangent



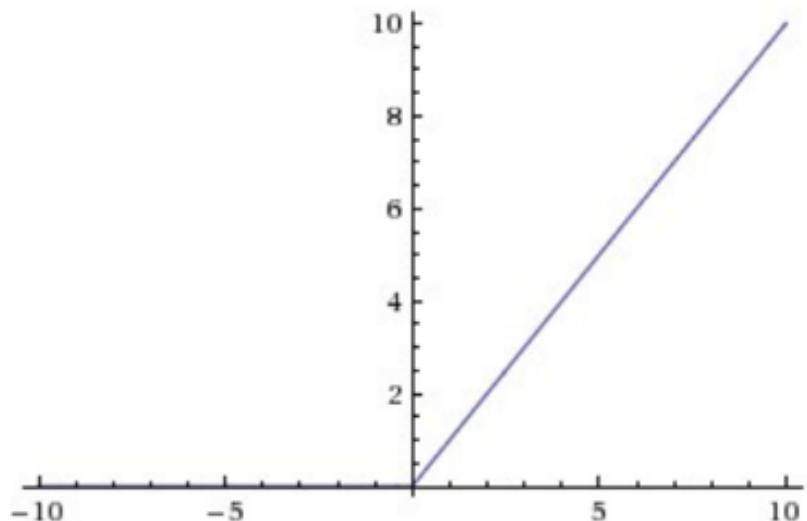
tanh(x)

- ▶ Squashes numbers to range [-1,1]
 - ▶ Zero-centered is good
- ▶ Still have practical issues

Activation functions

Rectified Linear Units

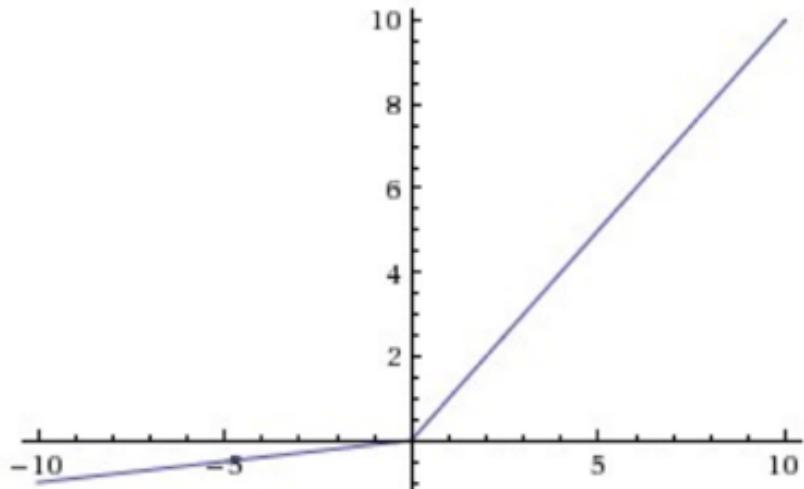
$$f(x) = \max(0, x)$$



- ▶ Most popular in current DL work
- ▶ Computationally efficient, ReLU neurons can be trained much faster
- ▶ Addresses the practical issues of sigmoid and tanh
- ▶ Neurons can “die” during training

Activation functions

Leaky ReLU



- ▶ Computationally efficient
- ▶ Will not “die” like ReLU

Leaky ReLU

Activation functions

In practice

Use ReLU

Try out Leaky ReLU

You can try out tanh but it probably won't work as well as ReLU

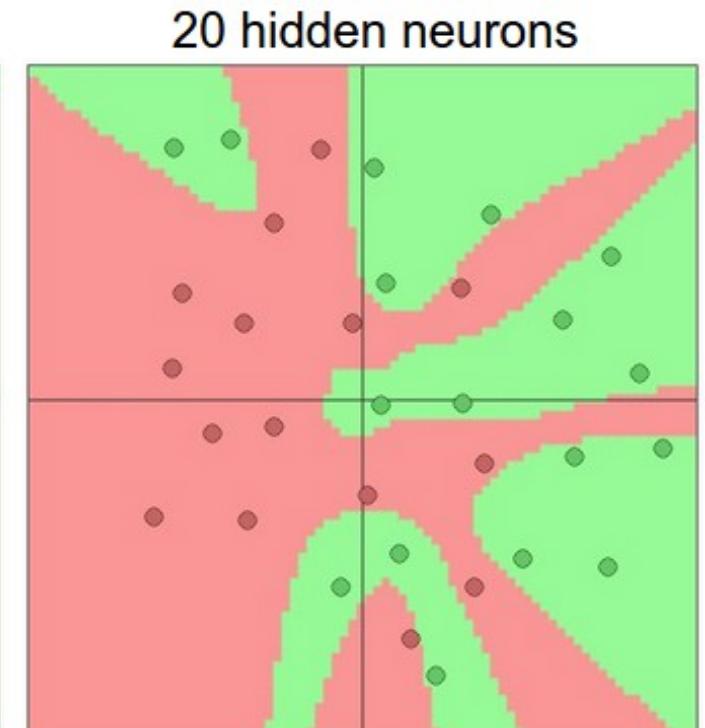
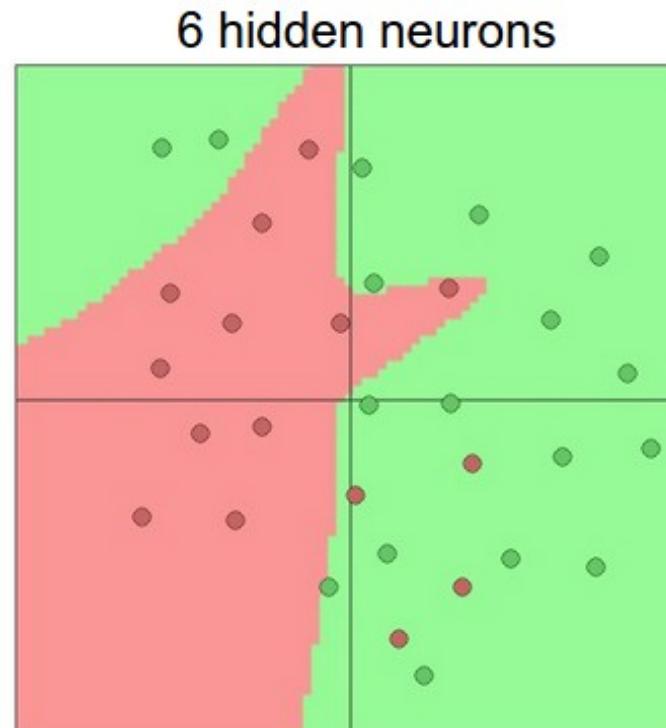
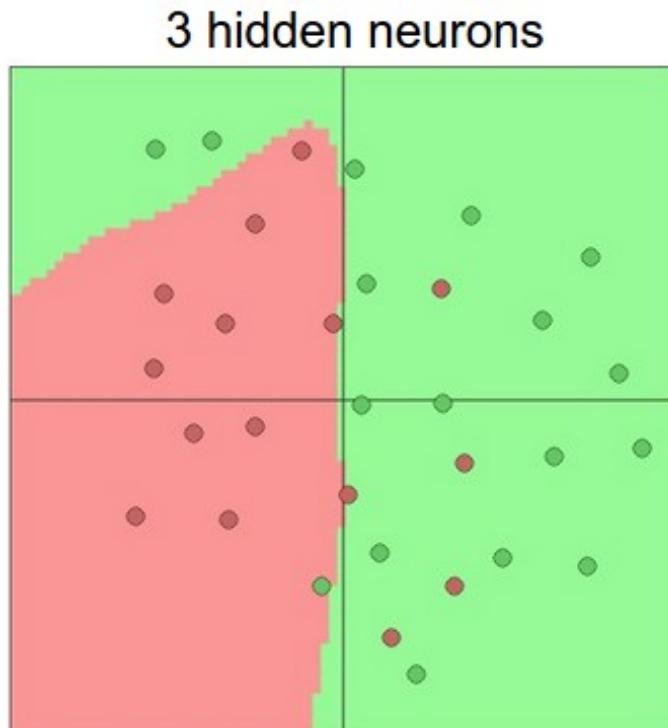
Never use sigmoid

Training Neural Networks

- Procedure for training Neural Networks
 - Perform inference on the training set
 - Calculate the error between the predictions and actual labels of the training set
 - Determine the contribution of each Neuron to the error
 - Modify the weights of the Neural Network to minimize the error
- Error contributions are calculated using Backpropagation
- Error minimization is achieved with Gradient Descent

Setting the number or layers and their size

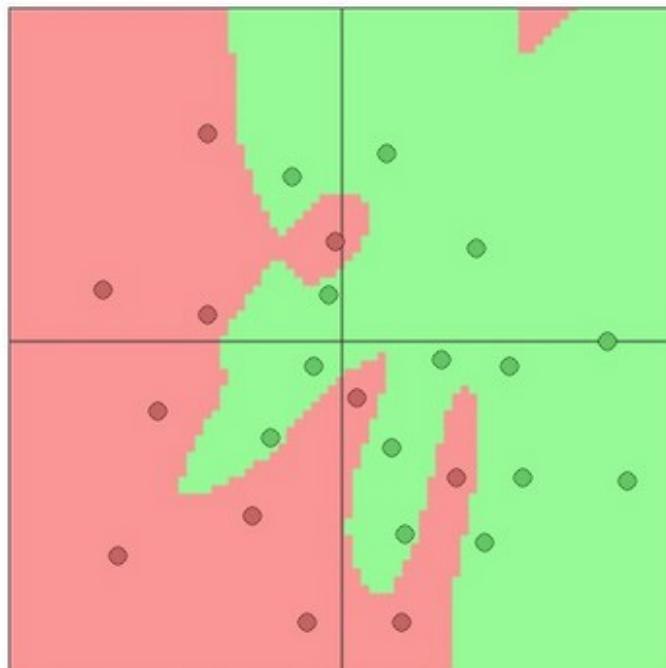
Overfitting, generalization and outliers



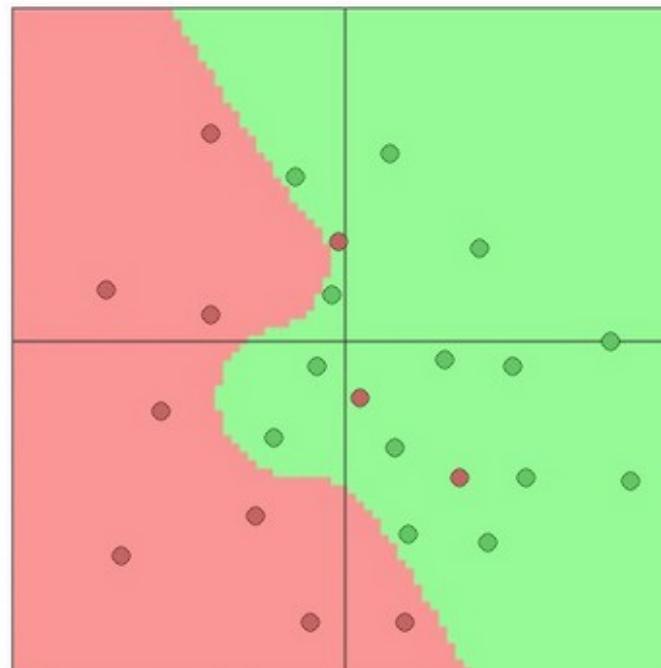
Effect of regularization

Overfitting, generalization and outliers

$\lambda = 0.001$



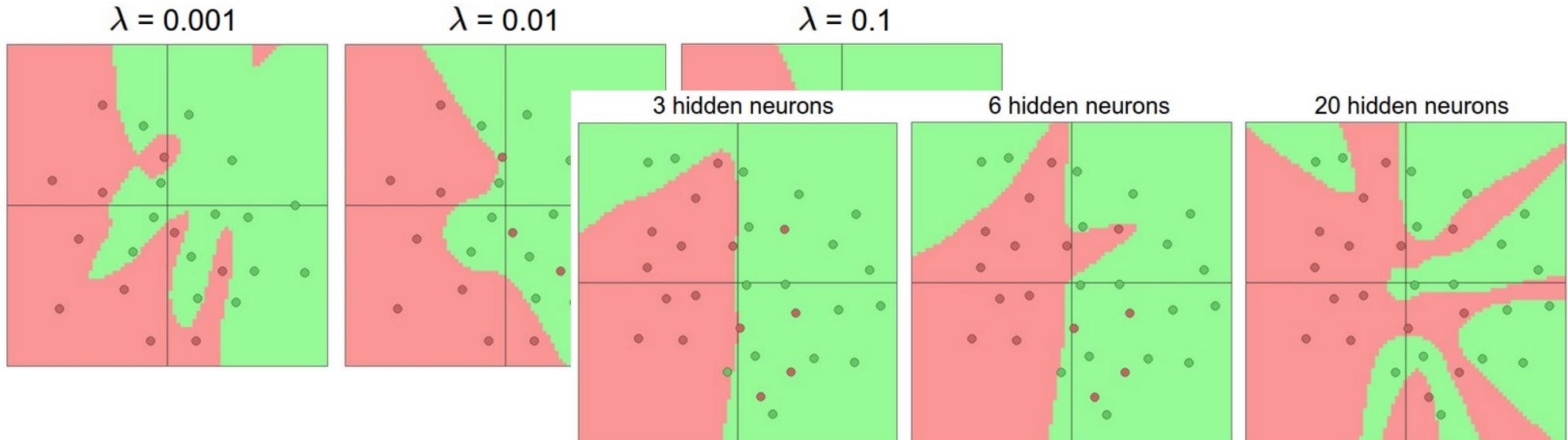
$\lambda = 0.01$



$\lambda = 0.1$



Network size and overfitting



Takeaway:

- Don't be afraid of overfitting: use a network as big as you can train
- Use regularization to control overfitting

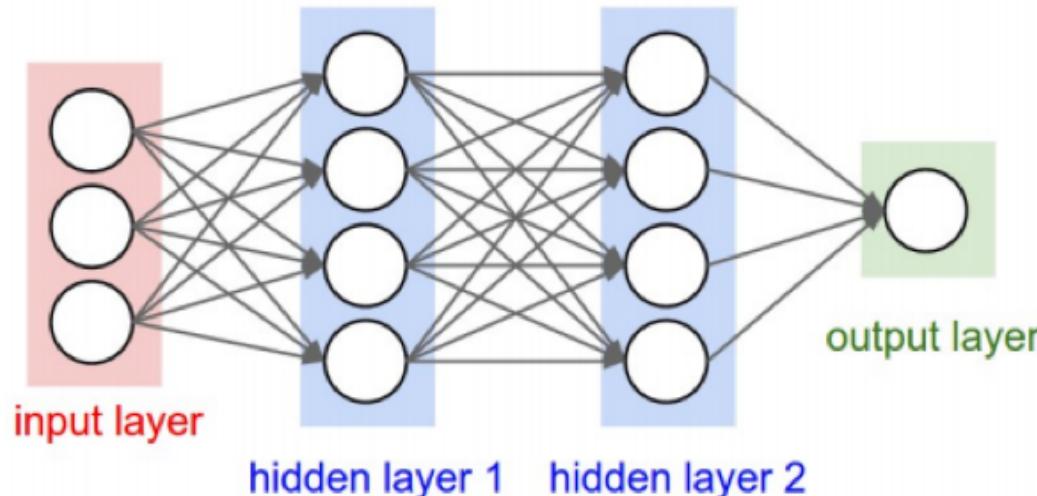
Vectorized feed-forward computation

```
1 f = lambda x: np.max(0,x)
2
3 class Neuron:
4
5     def neuron_fire(self, inputs):
6         """ assume inputs and weights are 1-D arrays and bias is a number """
7         cell_body_sum = np.dot(self.weights, inputs) + self.bias
8         activation = f(cell_body_sum) # ReLU activation function
9         return activation
```

Feed forward between two fully connected layers as a matrix-vector multiplication
(followed by activation)

This is very good for GPU acceleration

Example feed-forward computation of network

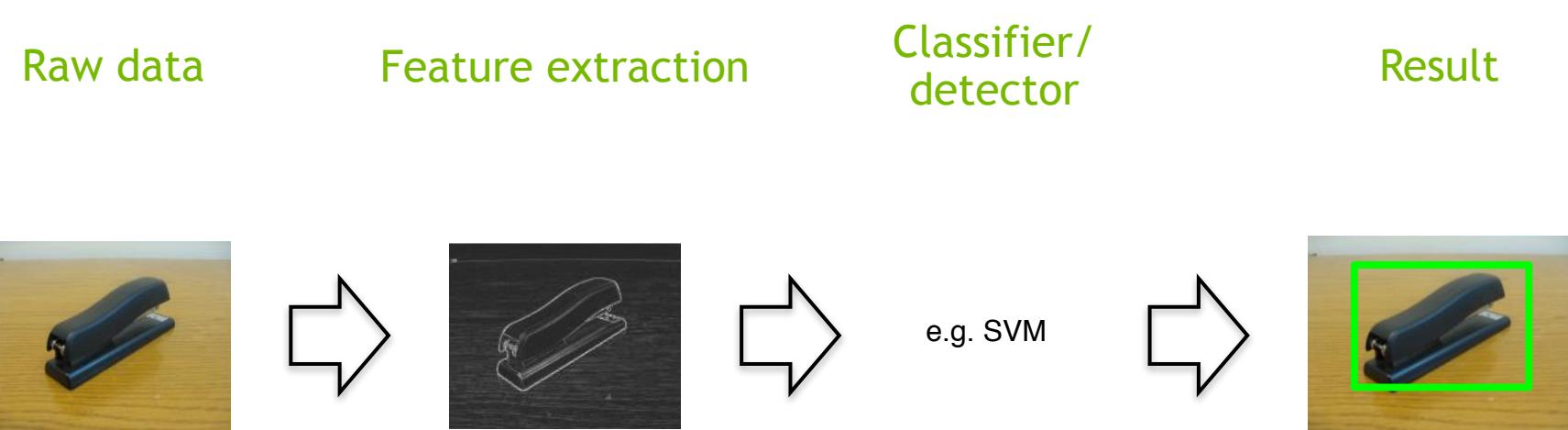


```
1 # forward-pass of a 3-layer neural network
2 f = lambda x: np.max(0,x)          # activation function (ReLU)
3 x = np.random.rand(3,1)            # random input vector of three numbers (3x1)
4 h1 = f(np.dot(W1, x) + b1)        # calculate first hidden layer activations (4x1)
5 h2 = f(np.dot(W2, h1) + b2)        # calculate second hidden layer activations (4x1)
6 out = f(np.dot(W3, h2) + b3)       # output neuron (1x1)
```

Conclusion

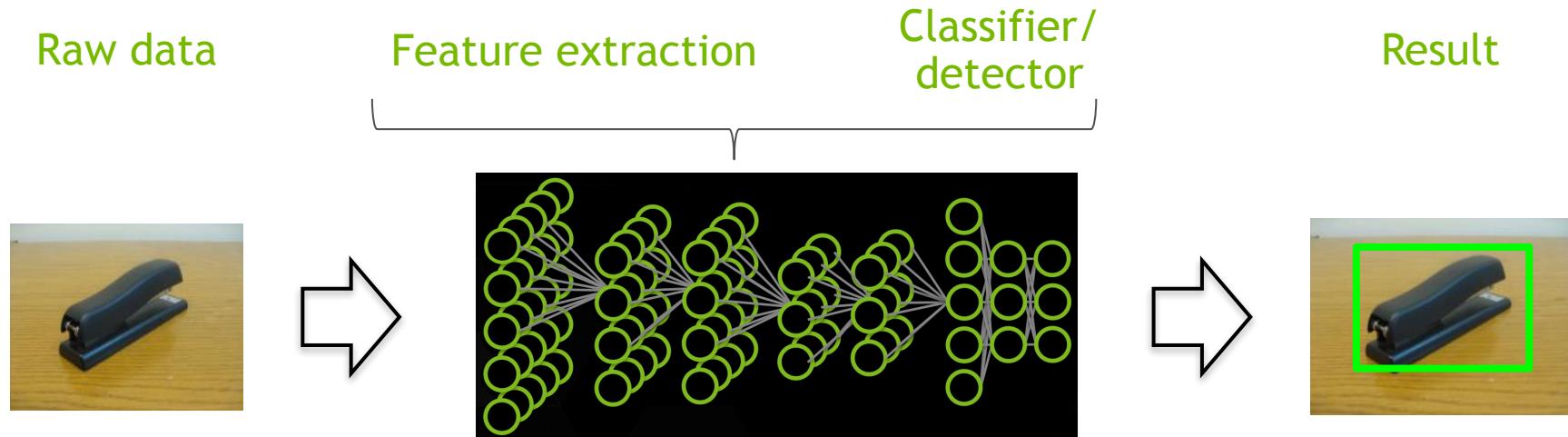
ANNs in Deep Learning

Traditional Machine Perception - hand tuned features



ANNs in Deep Learning

Deep Learning Machine Perception



Hyperparameter

You will see every point saying “it depends”

