

Neural Networks

ECE 4200

General steps we do

- $S = \{(\bar{X}_1, y_1), (\bar{X}_2, y_2), \dots, (\bar{X}_n, y_n)\}$
- Decide the algorithm/architecture to use (SVM/NN/LR).

Let \mathbf{W} be the algorithm's parameters (weights)

$A(\mathbf{W}, \bar{X})$: prediction of \mathbf{W} on \bar{X}

- Decide the loss function (L2, hinge loss, L2+Reg, binary cross entropy)

$$L(\mathbf{W}, (\bar{X}_i, y_i)) = L(A(\mathbf{W}, \bar{X}_i), y_i)$$

$$L(\mathbf{W}, S) = \sum_i L(\mathbf{W}, (\bar{X}_i, y_i))$$

- Find parameters that solve:

$$\min_{\mathbf{W}} L(\mathbf{W}, S)$$

Gradient descent

$$\min_{\mathbf{W}} L(\mathbf{W}, S)$$

This is hard to solve in most cases (no **closed form**).

$L(\mathbf{W}, S)$ is **just** a function of \mathbf{W}

In practice, use gradient descent (recall logistic regression).

- Start with some $\mathbf{W}(0)$
- Repeat until convergence:
 - Compute $\nabla L(\mathbf{W}, S)|_{\mathbf{w}(j)}$
 - $\mathbf{W}(j + 1) = \mathbf{W}(j) + \eta \cdot \nabla L(\mathbf{W}, S)|_{\mathbf{w}(j)}$

Pros, and cons.

Full gradient descent (GD)

Compute $\nabla L(\mathbf{W}, S)|_{\mathbf{W}(j)}$

$$\nabla L(\mathbf{W}, S) \Big|_{\mathbf{W}(j)} = \sum_i \nabla_{\mathbf{W}} L(\mathbf{W}, (\bar{X}_i, y_i)) \Big|_{\mathbf{W}(j)}$$

Each step evaluates the gradients for n points and **adds them up**

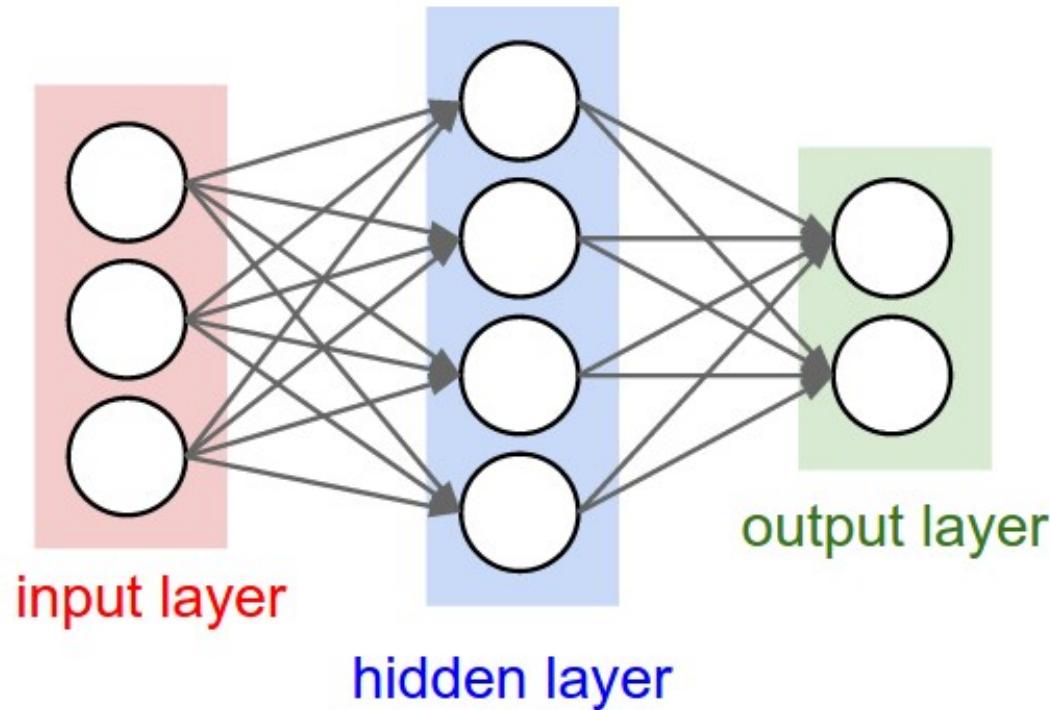
n can be very very large (Millions /Billions)

Full GD Prohibitive

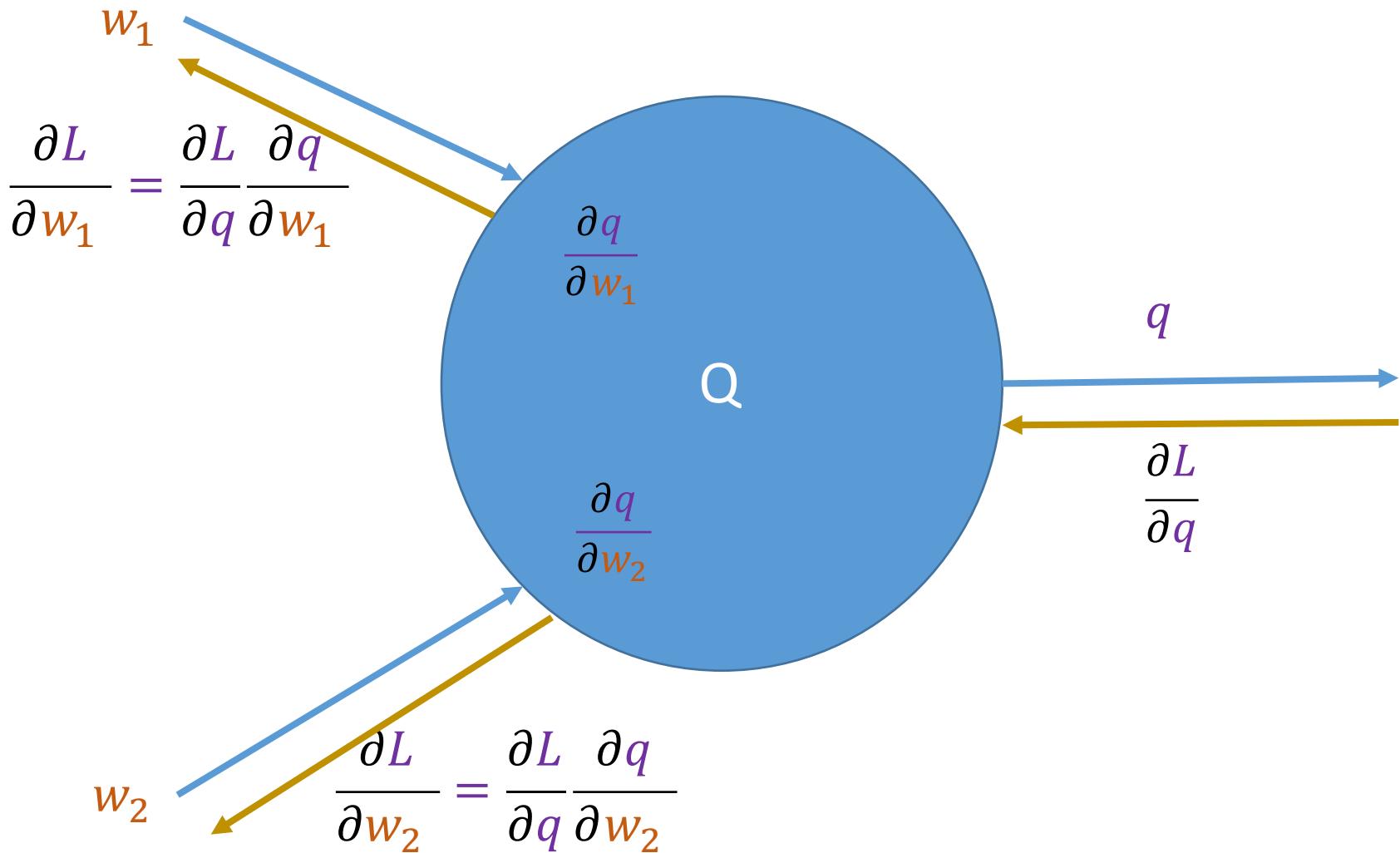
Neural Networks BP

Neural Networks

Large interconnected simple units



One Node



Sigmoid activation

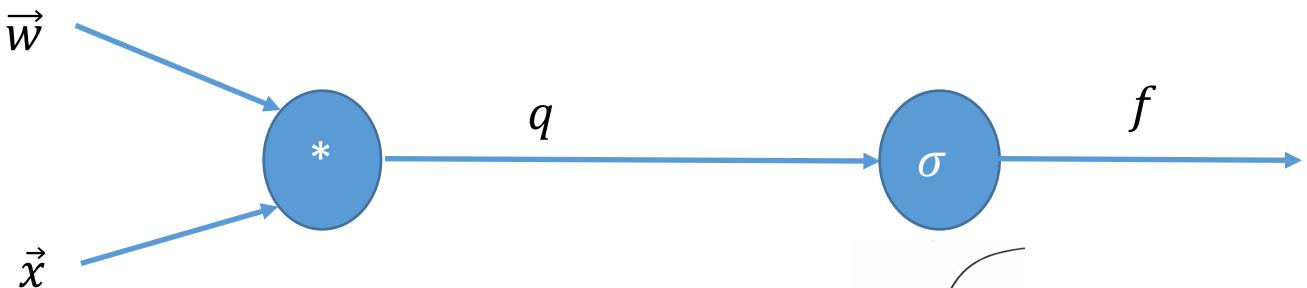
$$f(\bar{X}, \bar{w}) = \sigma(\bar{X} \cdot \bar{w}) = \frac{1}{1 + \exp(-(\bar{X} \cdot \bar{w}))}$$

One of the first and most popular activation functions.

Think of neurons:

Inputs are \vec{x} ,

Weights are \vec{w}

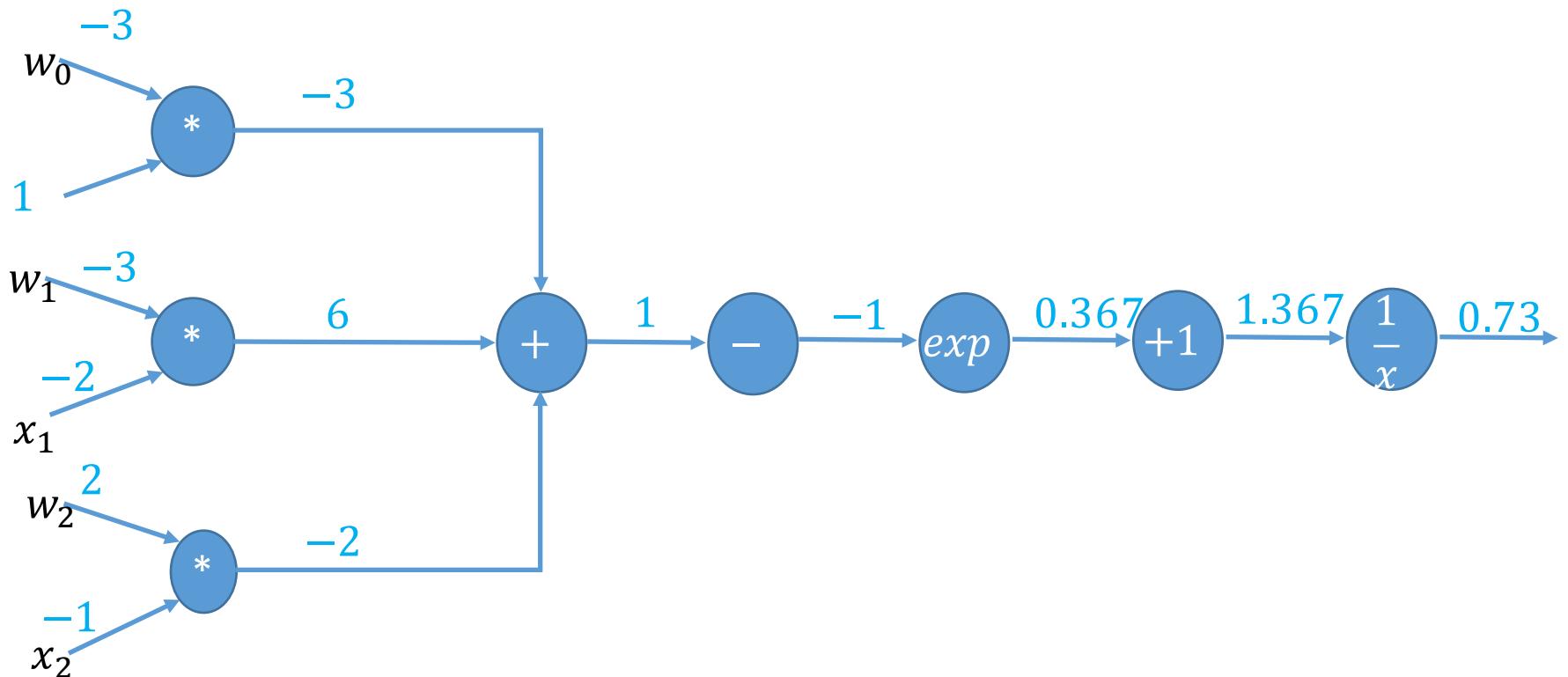


Linear combination of weights and inputs is activation fn
 $f(\vec{x}, \vec{w})$: the **probability** that the **neuron fires**

Sigmoid activation

$$d=3, \bar{w} = (w_0, w_1, w_2), \bar{X} = (1, x_1, x_2).$$

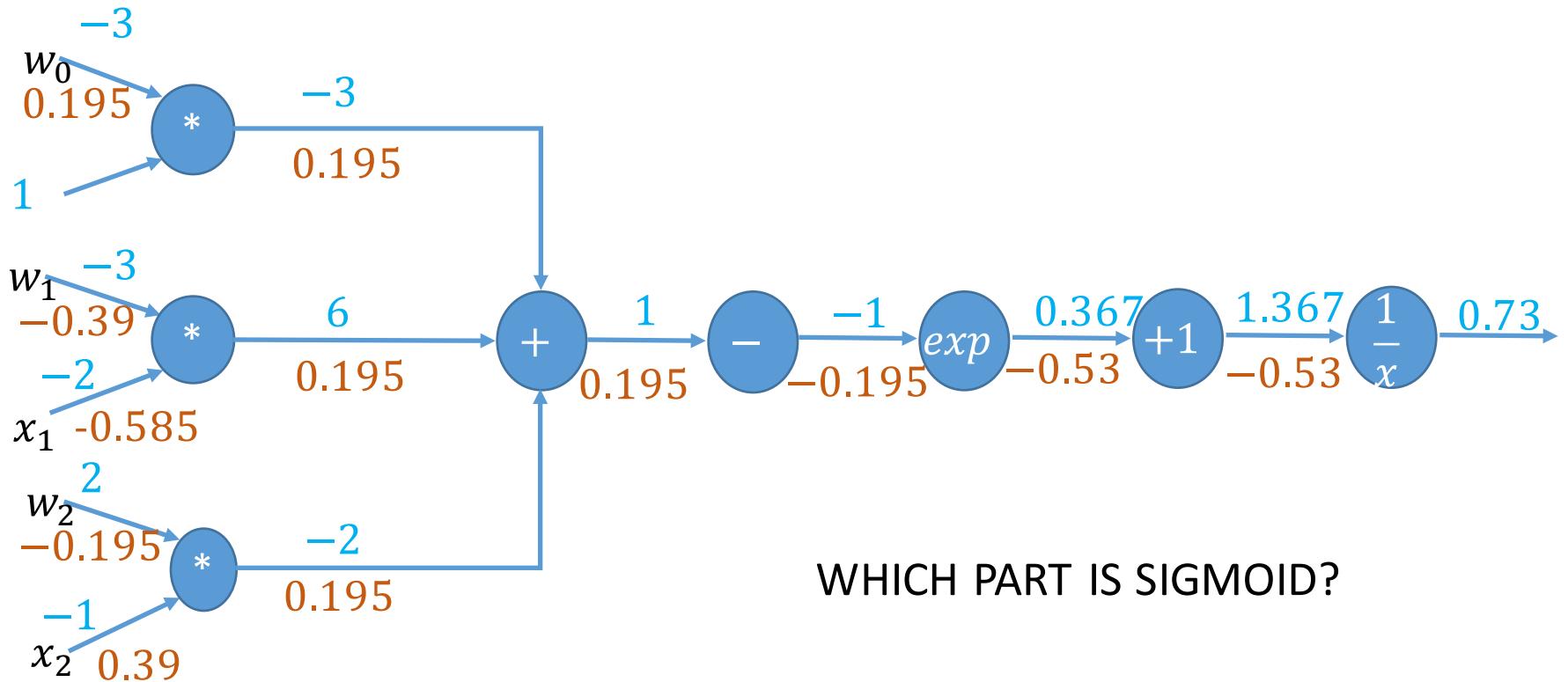
$$f(\vec{x}, \vec{w}) = \frac{1}{1 + \exp(-(w_0 + w_1 x_1 + w_2 x_2))}$$



Sigmoid activation

$$d=3, \bar{w} = (w_0, w_1, w_2), \bar{X} = (1, x_1, x_2).$$

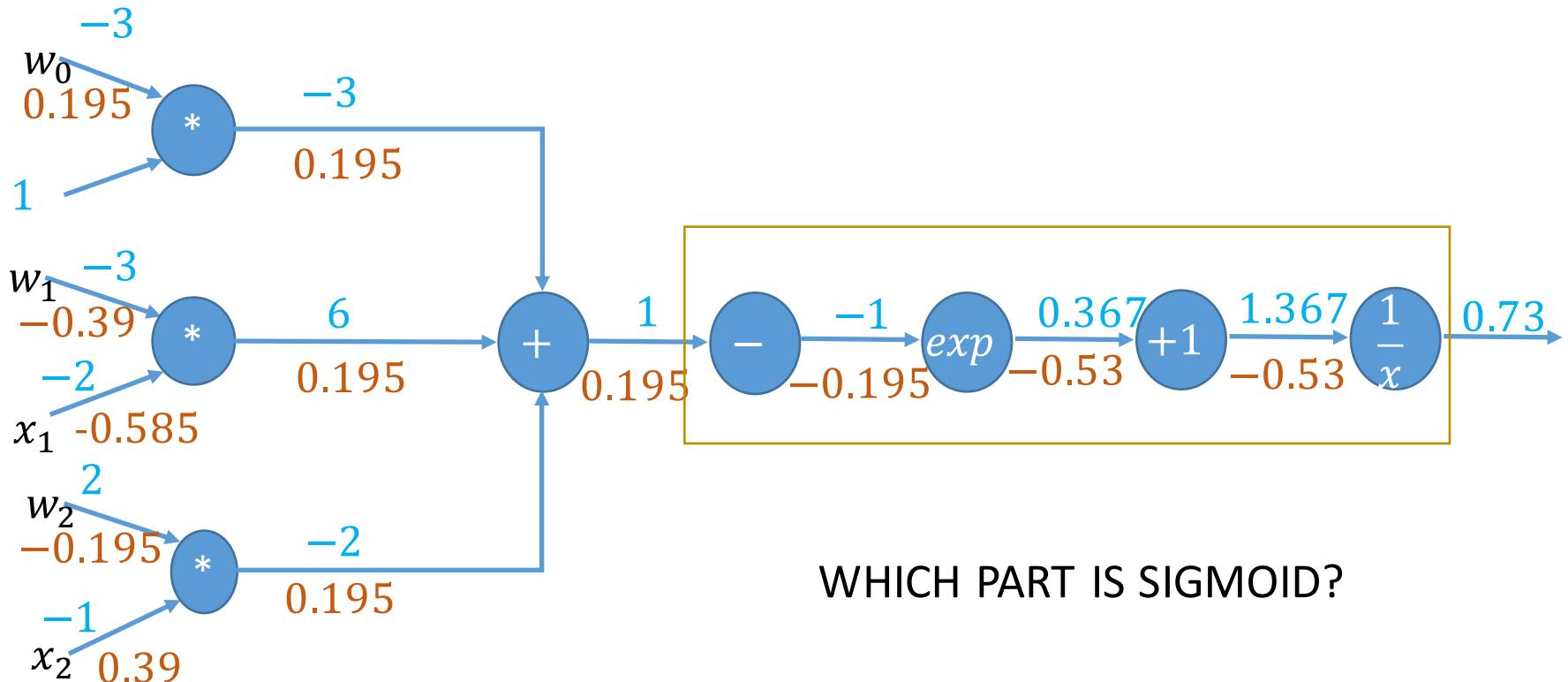
$$f(\bar{X}, \bar{w}) = \frac{1}{1 + \exp(-(w_0 + w_1 x_1 + w_2 x_2))}$$



Sigmoid activation

$$\bar{w} = (w_0, w_1, w_2), \bar{X} = (1, x_1, x_2).$$

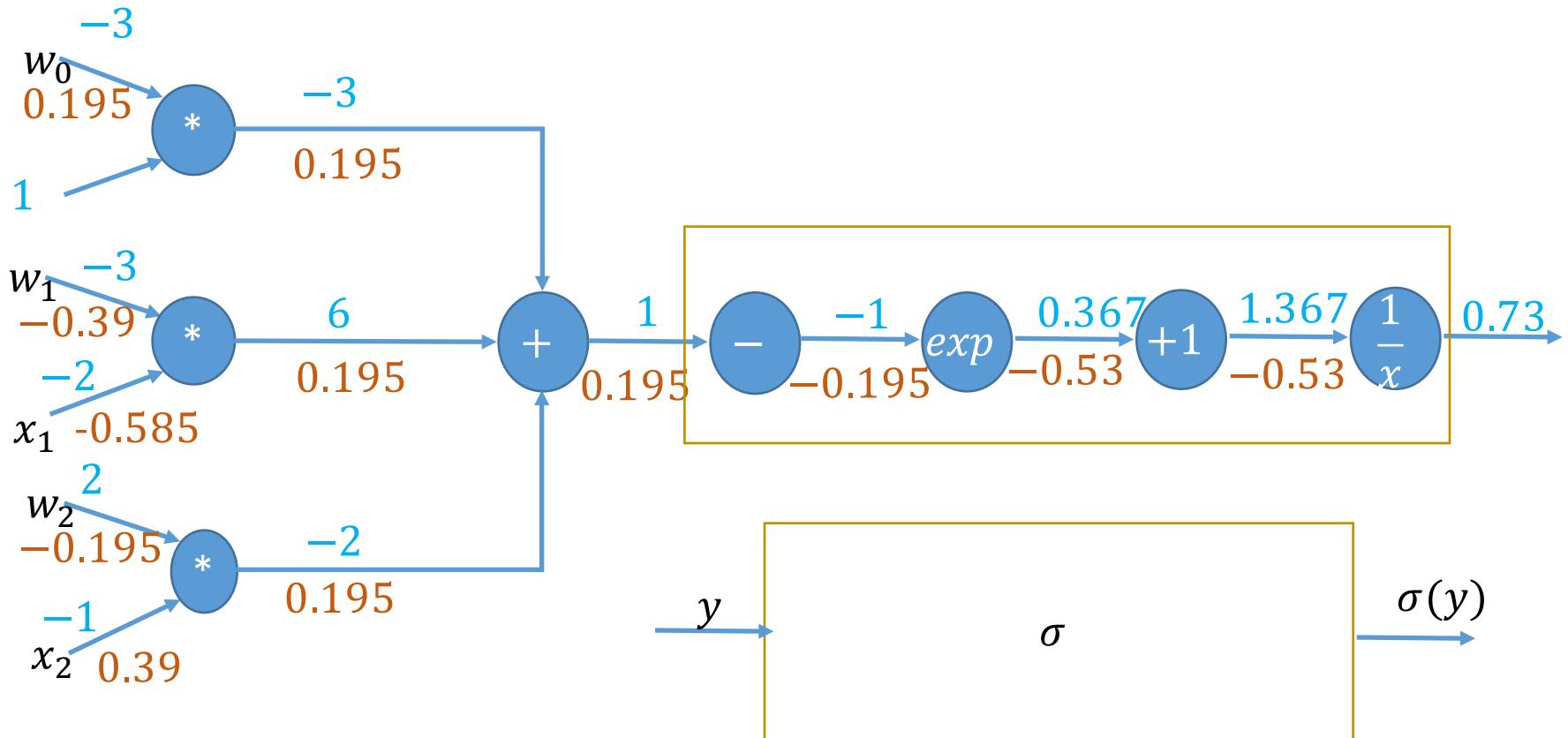
$$f(\bar{X}, \bar{w}) = \frac{1}{1 + \exp(-(w_0 + w_1 x_1 + w_2 x_2))}$$



Sigmoid activation

$$\bar{w} = (w_0, w_1, w_2), \bar{X} = (1, x_1, x_2).$$

$$f(\bar{X}, \bar{w}) = \frac{1}{1 + \exp(-(w_0 + w_1 x_1 + w_2 x_2))}$$



Derivative of sigmoid

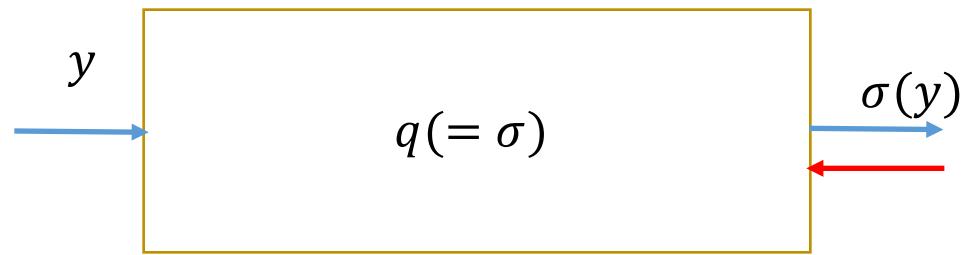
$$\sigma(y) = \frac{1}{(1 + e^{-y})}, 1 - \sigma(y) = \frac{e^{-y}}{(1 + e^{-y})}$$

$$\frac{d\sigma(y)}{dy} = \frac{d(1 + e^{-y})^{-1}}{dy} = \frac{-1}{(1 + e^{-y})^2} \frac{d(1 + e^{-y})}{dy} = \frac{e^{-y}}{(1 + e^{-y})^2}$$

$$\frac{d\sigma(y)}{dy} = \sigma(y) \cdot (1 - \sigma(y)).$$

Backpropagation for Sigmoid!

Sigmoid block



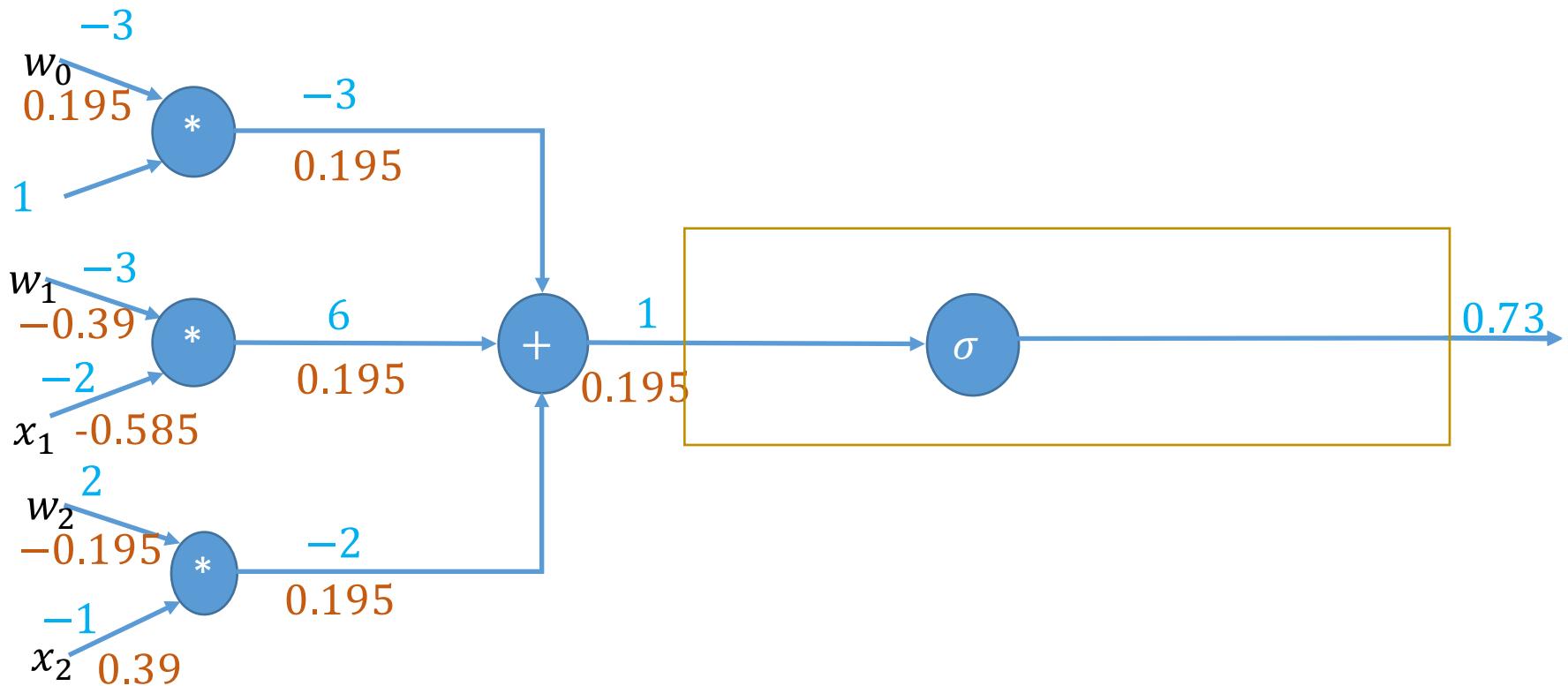
$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial q} \frac{\partial q}{\partial y}$$

$\frac{\partial L}{\partial q}$: comes from next layer

$$\frac{\partial q}{\partial y} = \sigma(y) \cdot (1 - \sigma(y))$$

Sigmoid activation

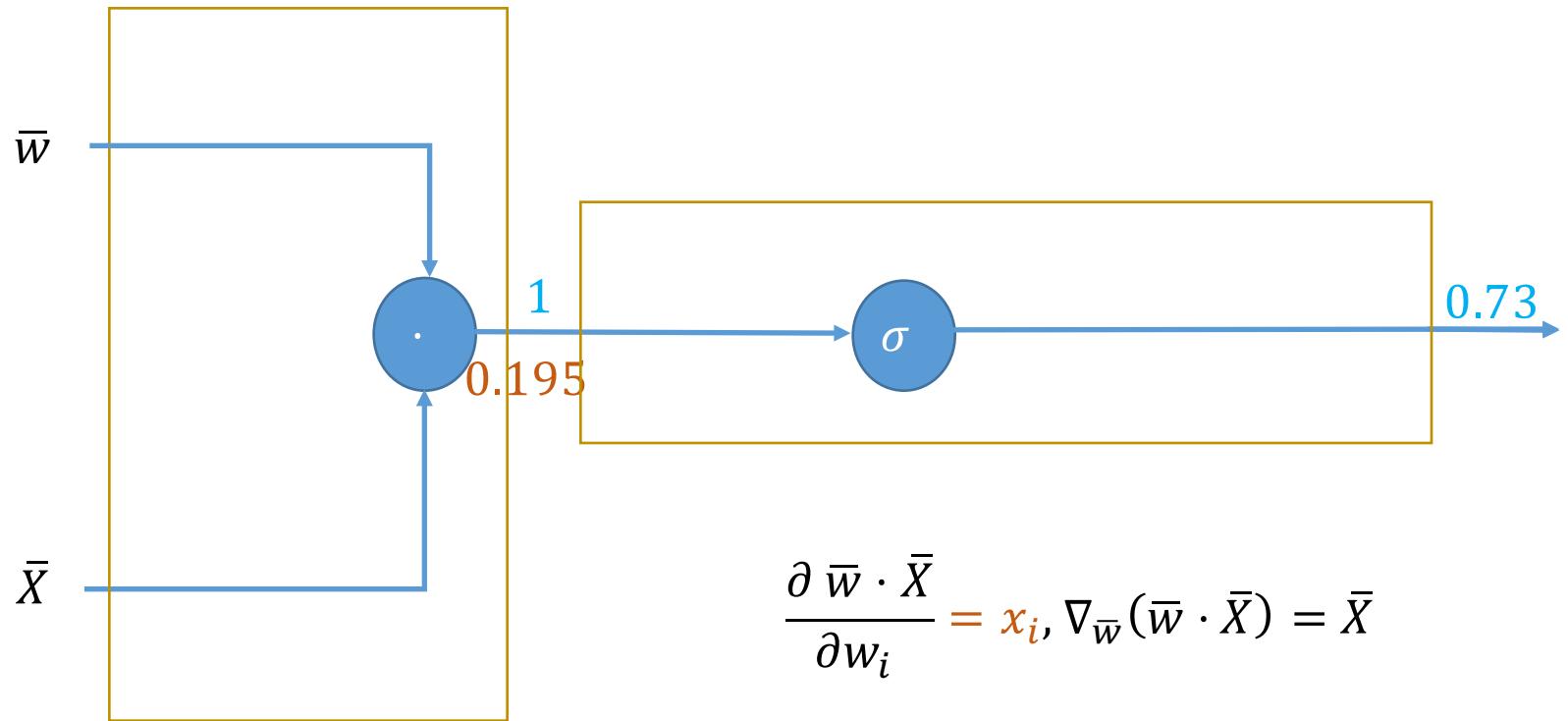
$$\bar{w} = (w_0, w_1, w_2), \bar{X} = (1, x_1, x_2).$$
$$f(\bar{X}, \bar{w}) = \sigma(w_0 + w_1 x_1 + w_2 x_2)$$



Sigmoid activation

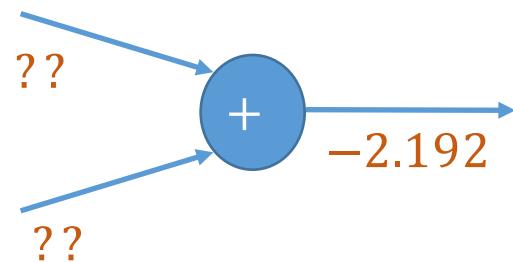
$$\bar{w} = (w_0, w_1, w_2), \bar{X} = (1, x_1, x_2).$$

$$f(\bar{X}, \bar{w}) = \sigma(\bar{w} \cdot \bar{X})$$



Derivatives for different *gates*

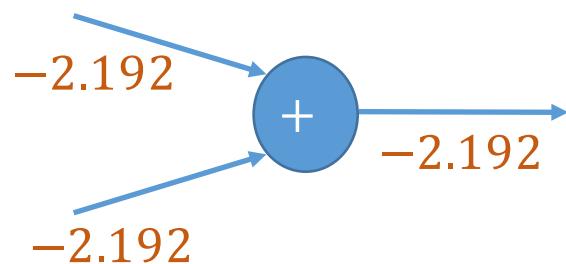
Add gate:



$$\frac{\partial \textcolor{violet}{L}}{\partial \textcolor{brown}{w}} = \frac{\partial \textcolor{violet}{L}}{\partial \textcolor{violet}{q}} \frac{\partial \textcolor{violet}{q}}{\partial \textcolor{brown}{w}}$$

Derivatives for different gates

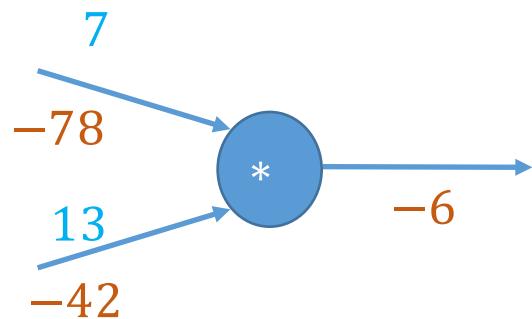
Add gate: gradient distributor



Derivatives for different gates

Add gate: gradient distributor

Mul gate: gradient switcher

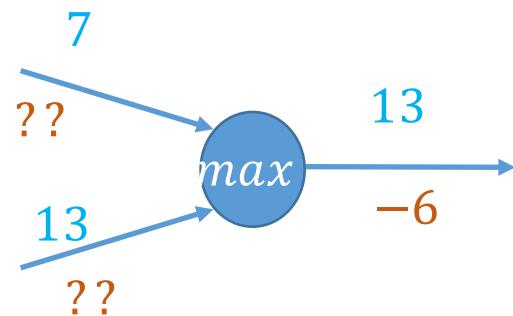


Derivatives for different gates

Add gate: gradient distributor

Mul gate: gradient switcher

Max gate:

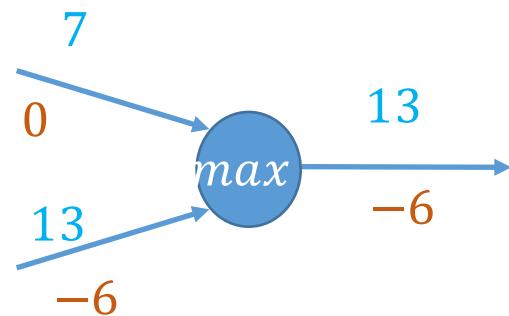


Derivatives for different gates

Add gate: gradient distributor

Mul gate: gradient switcher

Max gate:



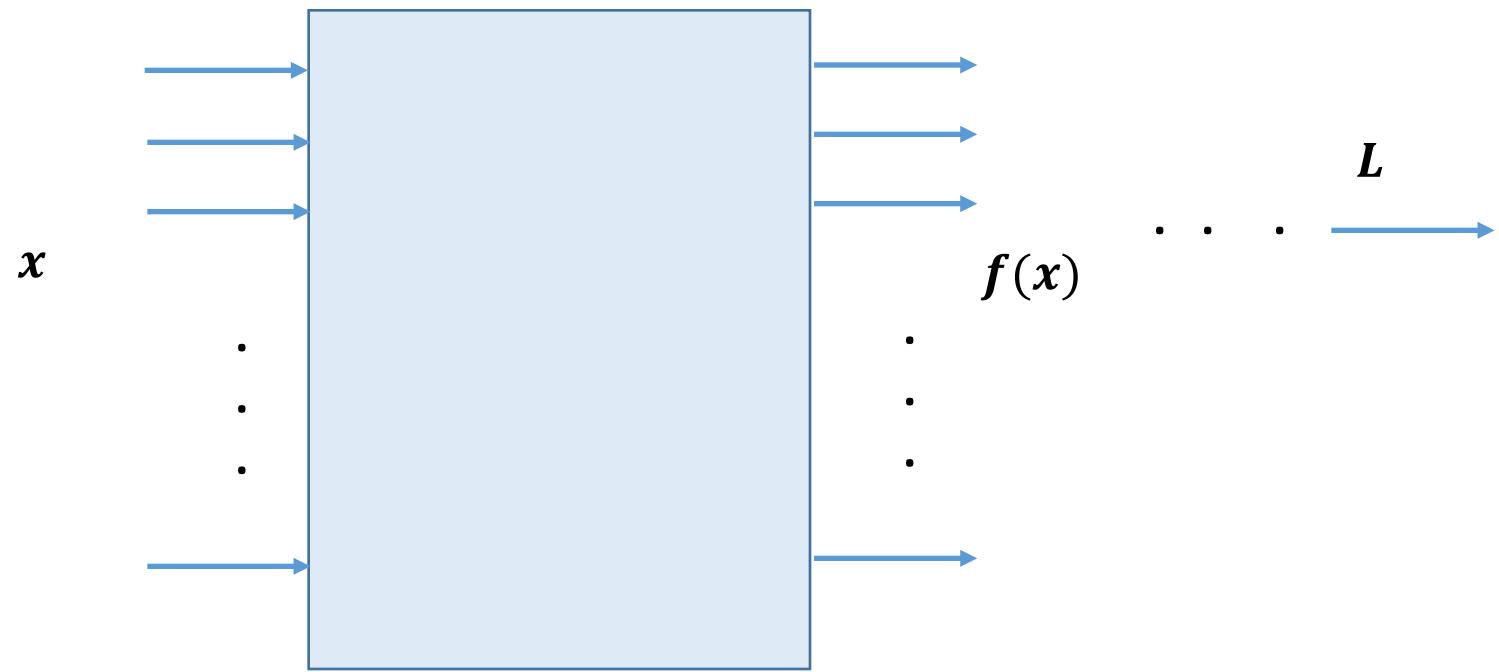
Plan for today

- Vector Outputs (Jacobian)
- Activation Functions
- Training NNs: Stochastic gradient descent, mini batch training
- Pre-processing
- CNN (just a primer)
- Loss Functions

Vector Outputs

Vector Outputs

Consider an operation/function of the form:



Functions $f: \mathbb{R}^d \rightarrow \mathbb{R}^m$, and L is real valued loss

Vector Outputs

Until now, $m = 1$, namely f was a scalar.

Chain rule:

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

The left hand side is actually the gradient, namely $\nabla_x L$.

How to compute **gradient** $\nabla_x L$ when f is vector valued?

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

Vector Outputs

Until now, $m = 1$, namely f was a scalar.

Chain rule:

$$\frac{\partial L}{\partial \mathbf{x}} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial \mathbf{x}}$$

The left hand side is actually the gradient, namely $\nabla_{\mathbf{x}} L$.

How to compute **gradient** $\nabla_{\mathbf{x}} L$ when $f \in \mathbb{R}^m$ is vector valued?

chain rule:

$$\frac{\partial L}{\partial \mathbf{x}} = \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \frac{\partial L}{\partial \mathbf{f}}$$

What are dimensions of these objects?

Jacobian

$\frac{\partial L}{\partial f}$: Gradient w.r.t f

Dimension of $\frac{\partial L}{\partial f}$?

$m \times 1$

$\frac{\partial f}{\partial x}$: Jacobian matrix

Dimension of $\frac{\partial f}{\partial x}$?

$d \times m$

Jacobian Matrix

$$\frac{\partial \mathbf{f}}{\partial \mathbf{x}} := \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_1}{\partial x_d} & \dots & \frac{\partial f_m}{\partial x_d} \end{bmatrix}$$

Dimension of $\frac{\partial \mathbf{f}}{\partial \mathbf{x}}$?

$d \times m$

Jacobian Matrix

You can also think in terms of the following equation:

$$\frac{\partial L}{\partial x_1} = \sum_{i=1}^m \frac{\partial L}{\partial f_i} \cdot \frac{\partial f_i}{\partial x_1}$$

Each entry comes from one entry in the matrix vector product.

Jacobian Matrix

$$\begin{bmatrix} \frac{\partial L}{\partial x_1} \\ \frac{\partial L}{\partial x_2} \\ \vdots \\ \frac{\partial L}{\partial x_d} \end{bmatrix} := \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_1}{\partial x_d} & \cdots & \frac{\partial f_m}{\partial x_d} \end{bmatrix} \begin{bmatrix} \frac{\partial L}{\partial f_1} \\ \frac{\partial L}{\partial f_2} \\ \vdots \\ \frac{\partial L}{\partial f_m} \end{bmatrix}$$

Dimensionality and computation?

Hopefully some structure present

Jacobian for ReLU

Example.

$f: \mathbb{R}^d \rightarrow \mathbb{R}^d$, element wise thresholding.

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

What is the complexity/dimensionality of: $\frac{\partial f}{\partial x}$

What is $\frac{\partial f}{\partial x}$?

This function forms a crucial part of **ReLU networks**

Linear Transformations

For a $m \times d$ matrix W , and $\mathbf{x} \in \mathbb{R}^d$

$$L(\mathbf{x}, W) = \| W \cdot \mathbf{x} \|_2^2$$

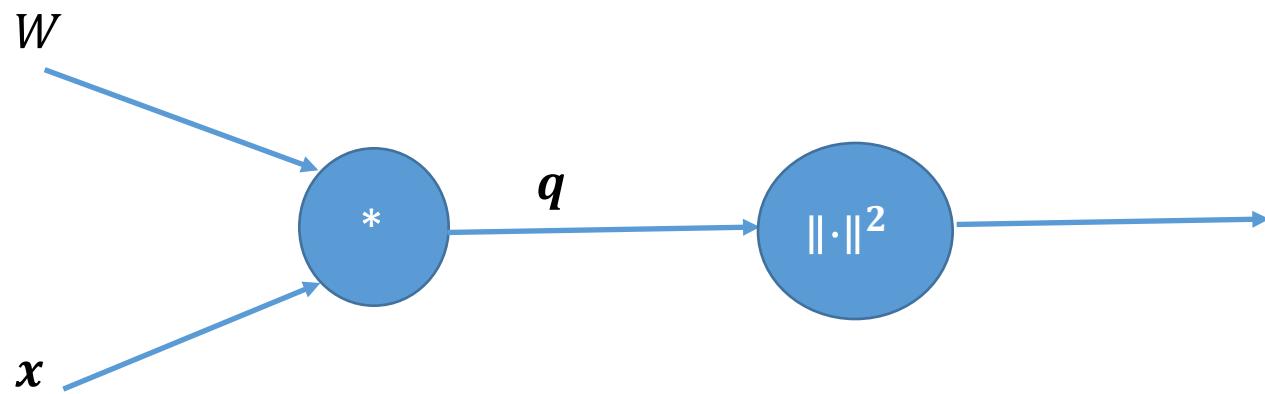
Computational graph for this?

Linear Transformations

For a $m \times d$ matrix W , and $x \in \mathbb{R}^d$

$$L(x, W) = \| W \cdot x \|_2^2$$

Computational graph for this?



Linear Transformations

$$L(x, W) = \|W \cdot x\|_2^2$$

$$W := \begin{bmatrix} W_{11} & \cdots & W_{1d} \\ \vdots & \ddots & \vdots \\ W_{m1} & \cdots & W_{md} \end{bmatrix}, \quad x = \begin{bmatrix} x_1 \\ \vdots \\ \vdots \\ x_d \end{bmatrix}$$

$$q = W \cdot x = \begin{bmatrix} W_{11}x_1 + \cdots + W_{1d}x_d \\ \vdots \\ \vdots \\ W_{m1}x_1 + \cdots + W_{md}x_d \end{bmatrix}$$

$$L(x, W) = \sum_i q_i^2 = q^T q = \sum_i (W_{i1}x_1 + \cdots + W_{id}x_d)^2$$

Computing gradients

$$L(x, W) = \sum_i q_i^2 = \mathbf{q}^T \mathbf{q} = \sum_i (W_{i1}x_1 + \cdots + W_{id}x_d)^2$$

$$\frac{\partial(q_1^2 + \cdots + q_m^2)}{\partial q_i} = 2q_i^2$$

$$\nabla_{\mathbf{q}} L(x, W) = \nabla_{\mathbf{q}} (\mathbf{q}^T \mathbf{q}) = 2\mathbf{q}.$$

Gradient dimension always matches what gradient is w.r.t.

Computing gradients

$$L(\mathbf{x}, W) = \sum_i \mathbf{q}_i^2 = \mathbf{q}^T \mathbf{q} = \sum_i (W_{i1}x_1 + \dots + W_{id}x_d)^2$$

$$\frac{\partial L(\mathbf{x}, W)}{\partial x_j} = \frac{\partial \sum_i \mathbf{q}_i^2}{\partial x_j} = \sum_i 2\mathbf{q}_i \frac{\partial \mathbf{q}_i}{\partial x_j} = \sum_i 2\mathbf{q}_i W_{ij}$$

$$\nabla_{\mathbf{x}} L(\mathbf{x}, W) = 2W^T \mathbf{q}$$

Do it differently: $L(\mathbf{x}, W) = \mathbf{q}^T \mathbf{q} = \mathbf{x}^T W^T W \mathbf{x}$

Gradient w.r.t. \mathbf{x} will be $2W^T W \mathbf{x} = 2W^T \mathbf{q}$

Computing gradients

$$L(\mathbf{x}, W) = \sum_i \mathbf{q}_i^2 = \mathbf{q}^T \mathbf{q} = \sum_i (W_{i1}x_1 + \cdots + W_{id}x_d)^2$$

$$\frac{\partial L(\mathbf{x}, W)}{\partial W_{jk}} = \frac{\partial \sum_i \mathbf{q}_i^2}{\partial W_{jk}} = \sum_i 2\mathbf{q}_i \frac{\partial \mathbf{q}_i}{\partial W_{jk}} = 2\mathbf{q}_j \mathbf{x}_k$$

$$\nabla_W L(\mathbf{x}, W) = 2\mathbf{q}\mathbf{x}^T.$$

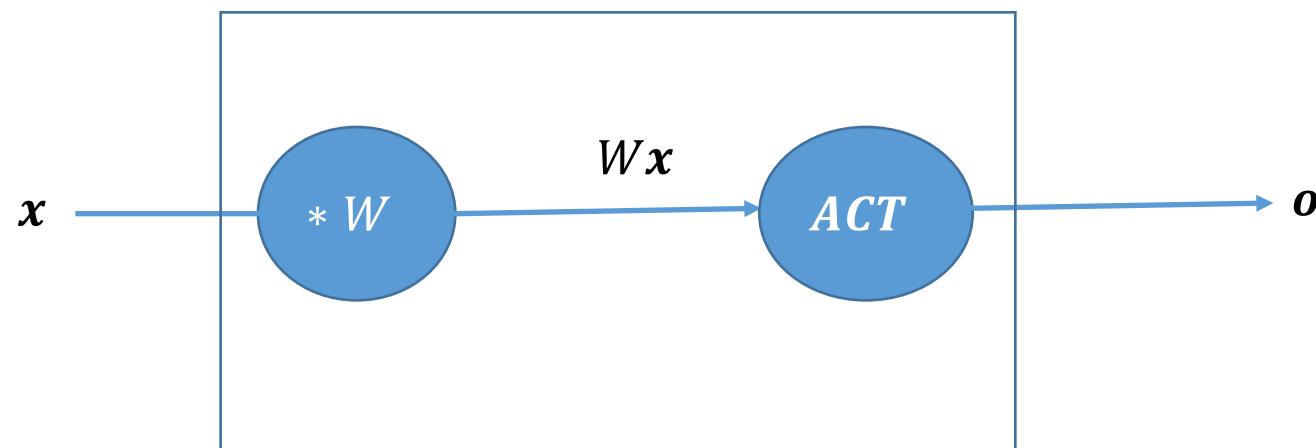
*Gradient dimension always matches what gradient is w.r.t.
Do it differently: $L(\mathbf{x}, W) = \mathbf{q}^T \mathbf{q} = \mathbf{x}^T W^T W \mathbf{x}$
Gradient w.r.t. W will be what ??*

Neural Networks

Linear transformations:

At a node, input x , and matrix W , output is Wx .

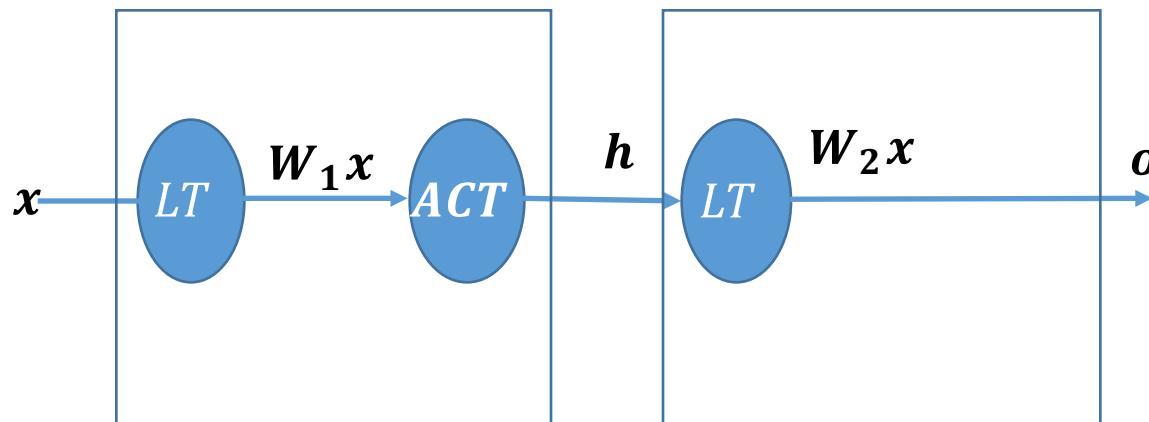
One LAYER:



Neural Networks

One LAYER:

Linear transformation, followed by non-linear transformation

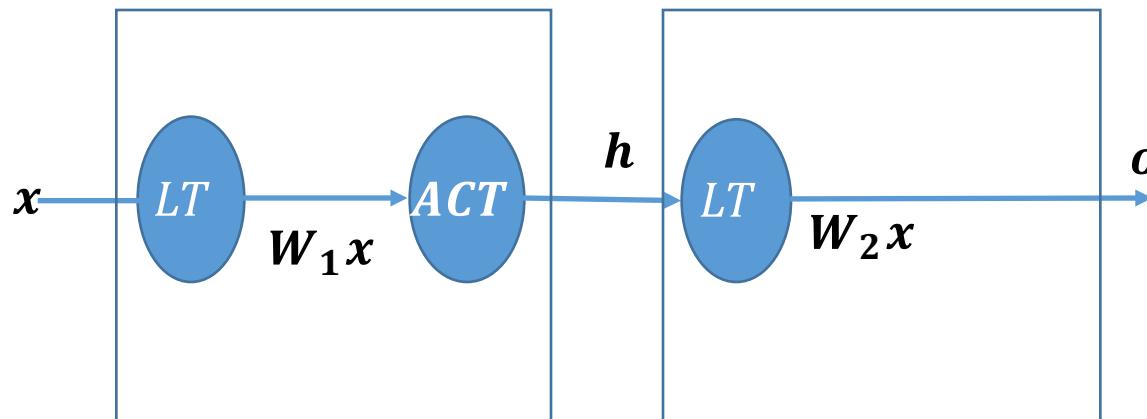


ReLU activation: $\mathbf{o} = \max \{0, W\mathbf{x}\}$, saves the non-negative terms

Neural Networks

One LAYER:

Linear transformation, followed by non-linear transformation

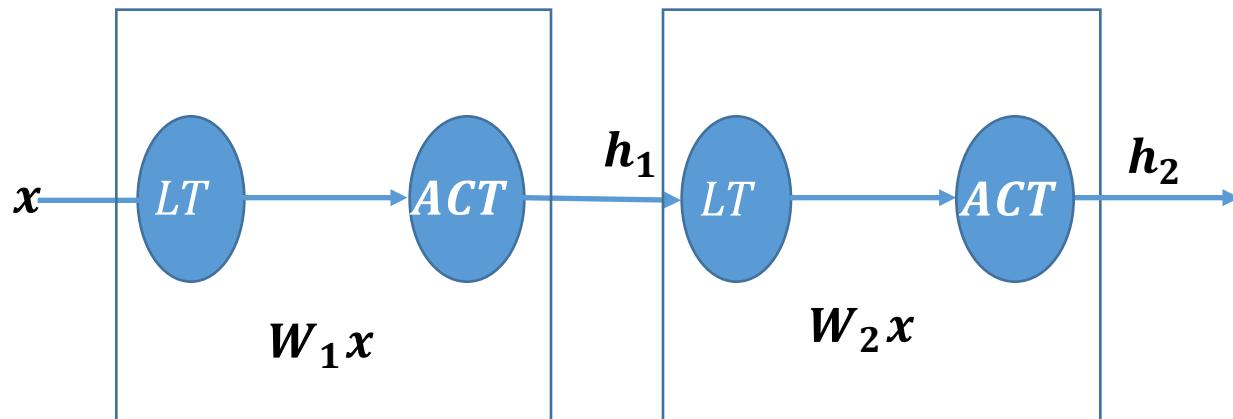


$$h = W_1x$$

2 layer network: $o = W_2 \cdot \max\{0, W_1x\}$, ReLU networks

Neural Networks

Can now cascade many layers ...



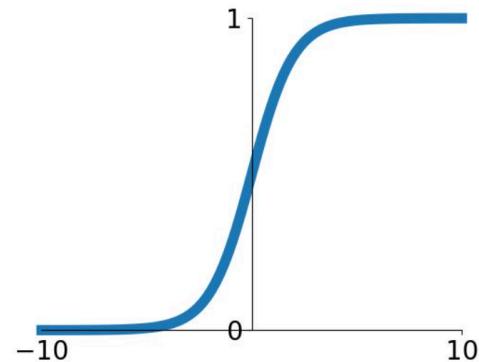
3 layer network: $o = W_3 \cdot \max\{0, W_2 \cdot \max\{0, W_1 x\}\}$, ReLU networks

Activation Functions

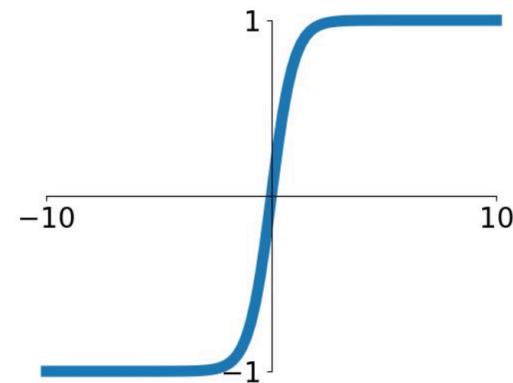
Activation functions/Non-linearity

What are the activations used?

Sigmoid: $\sigma(y) = \frac{1}{1+e^{-y}}$



Tanh: $\tanh y = \frac{e^y - e^{-y}}{e^y + e^{-y}}$



Pros and Cons

Saturated neurons kill gradients

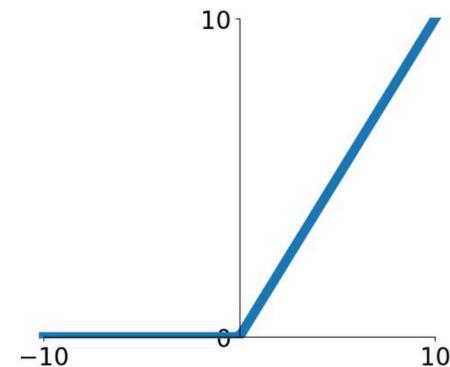
Sigmoid not super easy to compute

Not centered (why do we care?)

ReLU, Rectified Linear Units

What are the activations used?

$$\text{ReLU}(y) = \max(0, y)$$



Pros, and cons?

Does not saturate

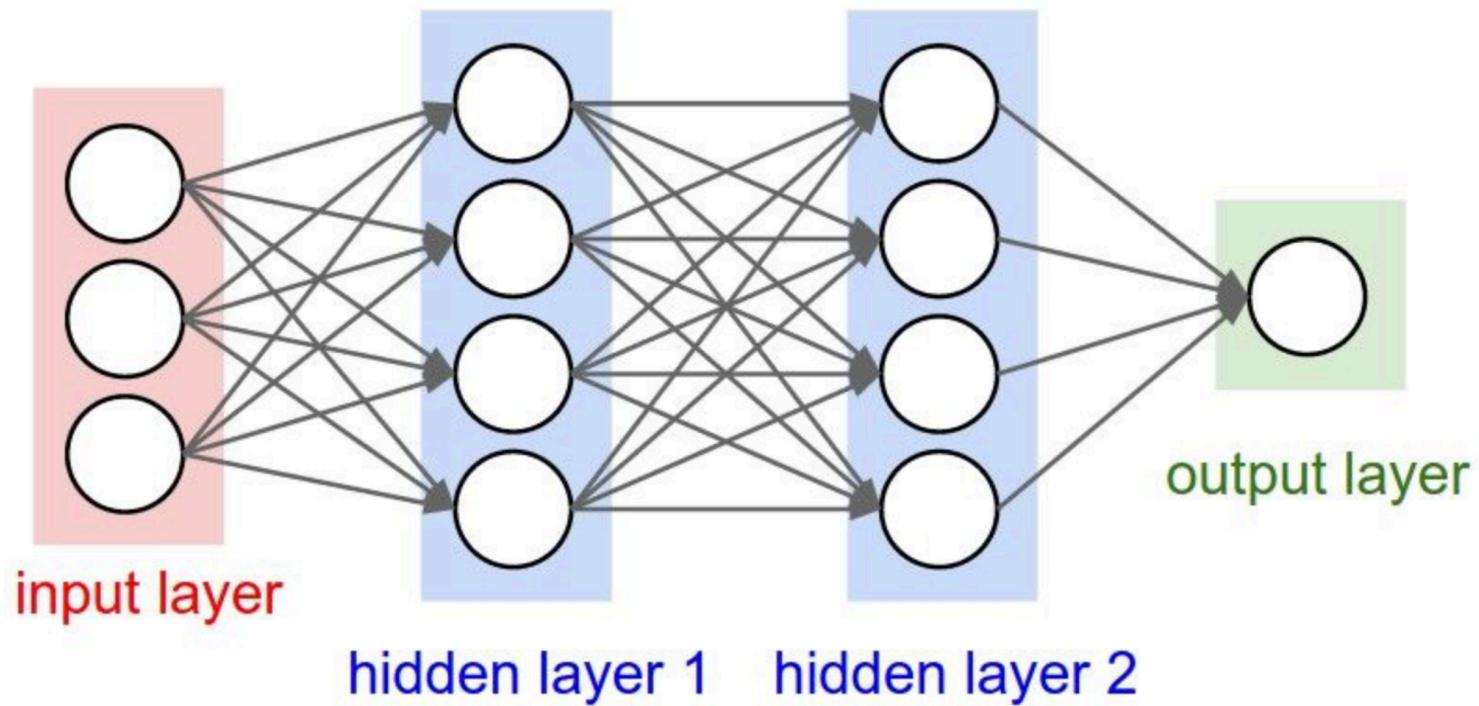
Fast to compute

Works great!

Many work **arounds**, no time to go **around** all of them

Training Neural Networks

3 Layer Neural Networks



Training Neural Networks

Fix the architecture (# Layers, size of each layer, loss functions, ...)

This tells you what the parameters are

Training NN using minibatch:

- Take a mini batch T
- Compute the functions in the computational graph, obtain loss
- Back propagate to calculate the gradients
- Update (each layer) based on the gradients

Repeat the process until convergence

Gradient descent

$$\min_{\mathbf{W}} L(\mathbf{W}, S)$$

This is hard to solve in most cases (no **closed form**).

$L(\mathbf{W}, S)$ is **just** a function of \mathbf{W}

In practice, use gradient descent (recall logistic regression).

- Start with some $\mathbf{W}(0)$
- Repeat until convergence:
 - Compute $\nabla L(\mathbf{W}, S)|_{\mathbf{W}(j)}$
 - $\mathbf{W}(j + 1) = \mathbf{W}(j) + \eta \cdot \nabla L(\mathbf{W}, S)|_{\mathbf{W}(j)}$

SGD

At iteration j:

Pick a random training (\bar{X}_i, y_i) , then

$$\mathbf{W}(j+1) = \mathbf{W}(j) + \eta \cdot \nabla L(\mathbf{W}, (\bar{X}_i, y_i)) \Big|_{\mathbf{W}(j)}$$

SGD + Momentum:

$$\mathbf{W}(j+1) = \mathbf{W}(j) + \eta \cdot \nabla L(\mathbf{W}, (\bar{X}_i, y_i)) \Big|_{\mathbf{W}(j)} + \gamma \cdot (\mathbf{W}(j) - \mathbf{W}(j-1))$$

Minibatch SGD

At each round:

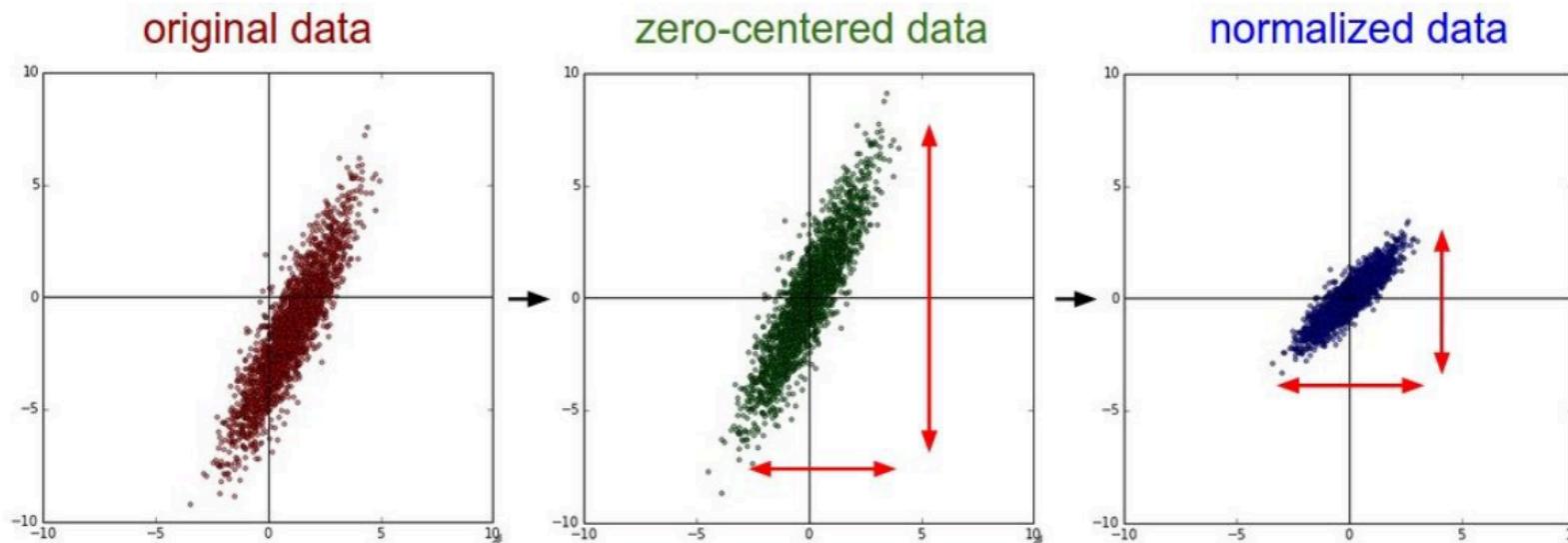
- Pick a few (tens/100's) training examples at random, say $T \subset S$
- Compute the gradient only for this batch:

$$\sum_{i \in T} \nabla_{\mathbf{W}} L(\mathbf{W}, (\bar{X}_i, y_i)) \Big|_{\mathbf{W}(j)}$$

- $\mathbf{W}(j+1) = \mathbf{W}(j) + \eta \cdot \nabla L(\mathbf{W}, T) \Big|_{\mathbf{W}(j)}$

Data Pre-processing

The first steps

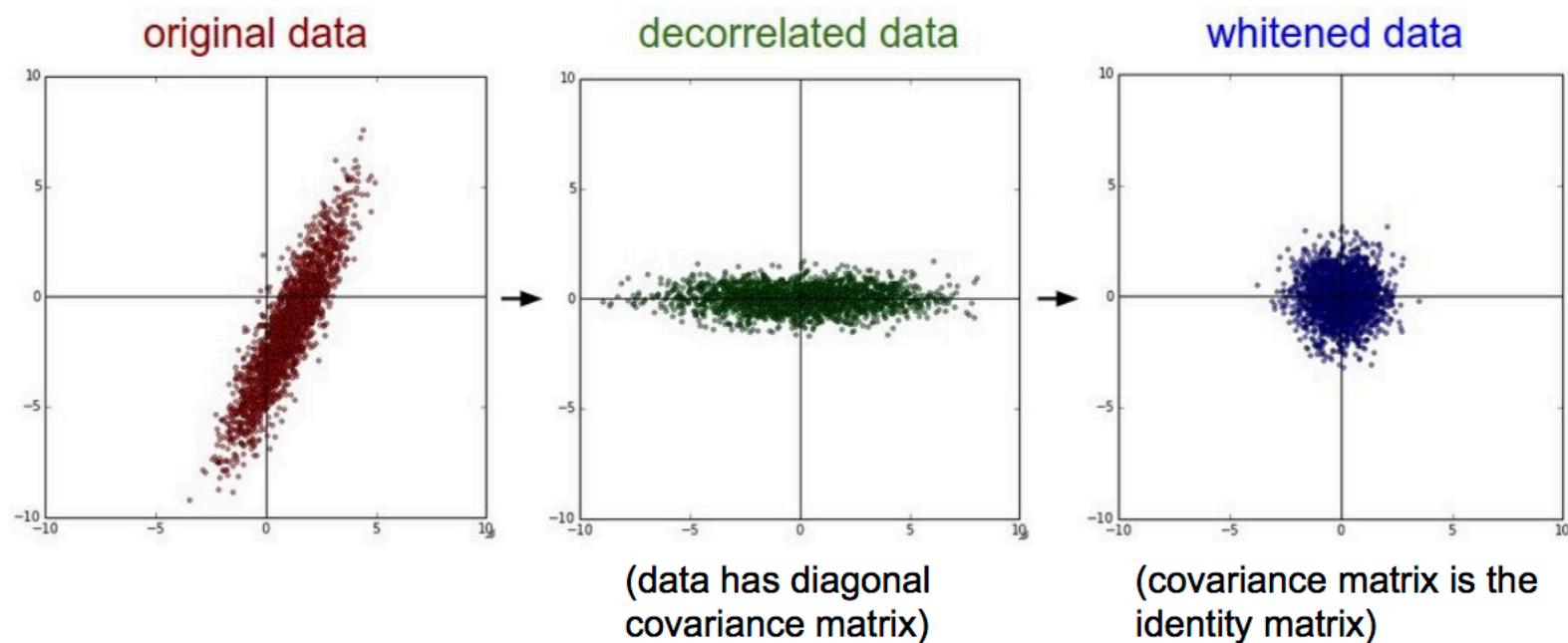


```
X -= np.mean(X, axis = 0)
```

```
X /= np.std(X, axis = 0)
```

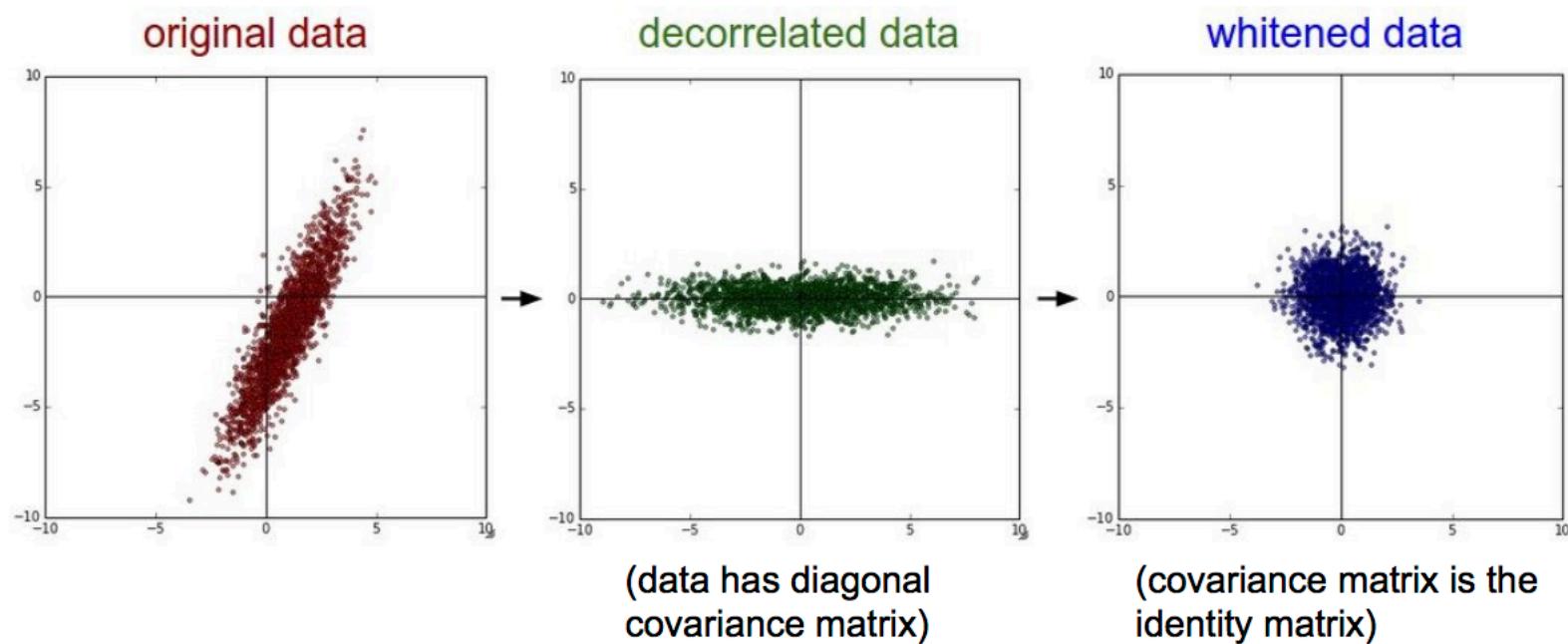
The first steps

In practice, you may also see **PCA** and **Whitening** of the data



The first steps

In practice, you may also see **PCA** and **Whitening** of the data



CNN

Next: Convolutional Neural Networks

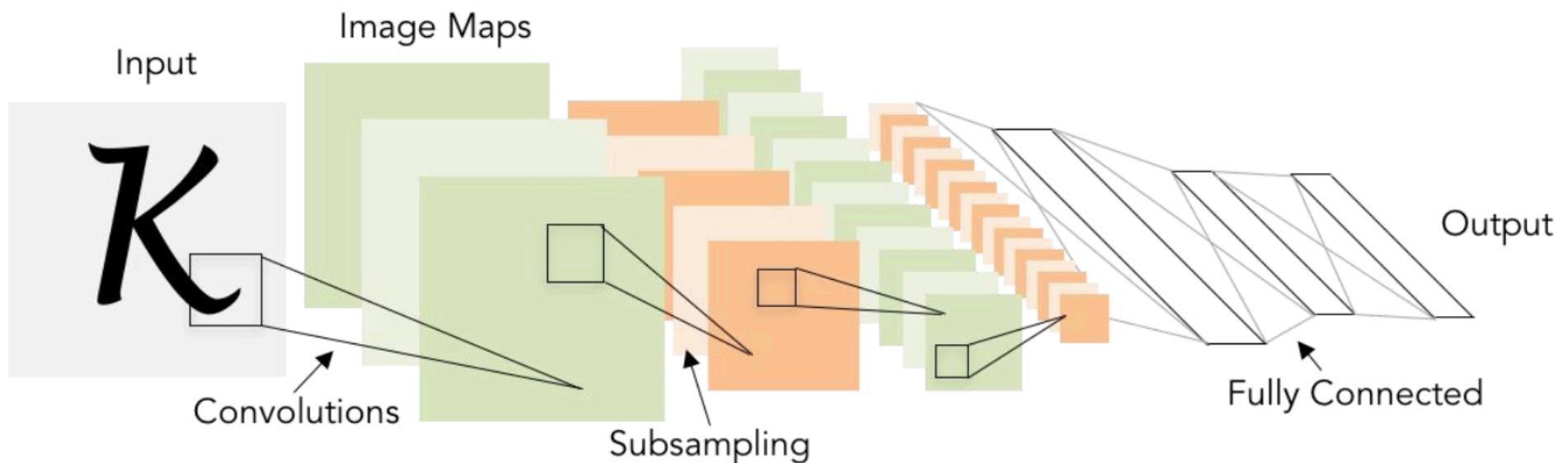


Illustration of LeCun et al. 1998 from CS231n 2017 Lecture 1

CNN

Not super new ... but super powerful now

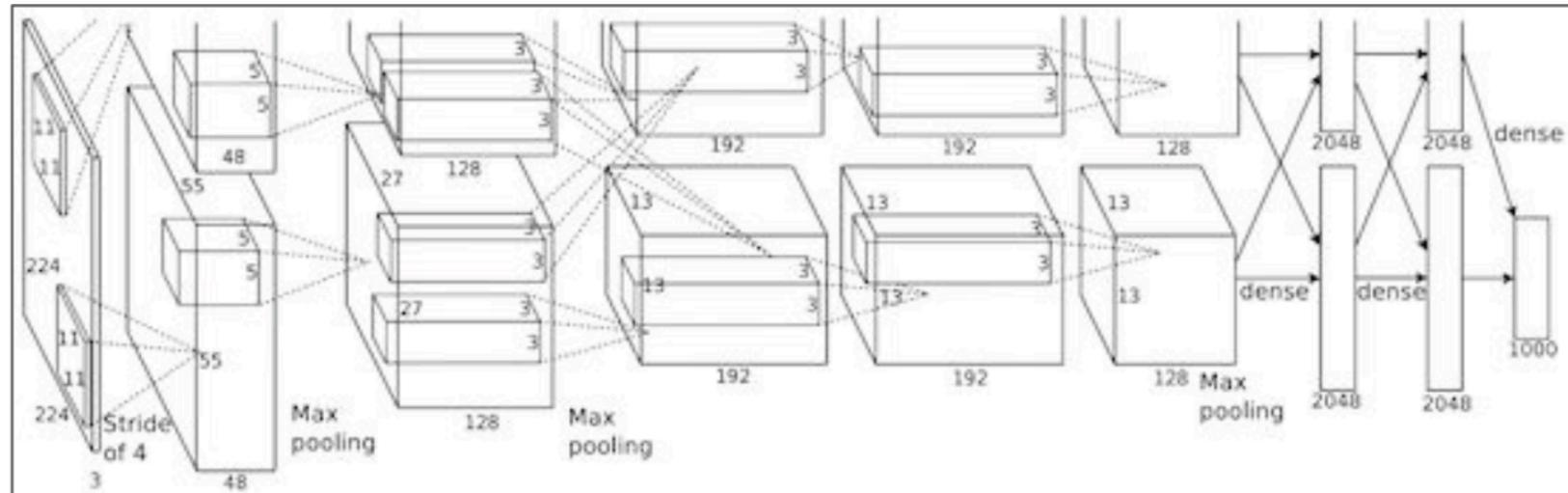
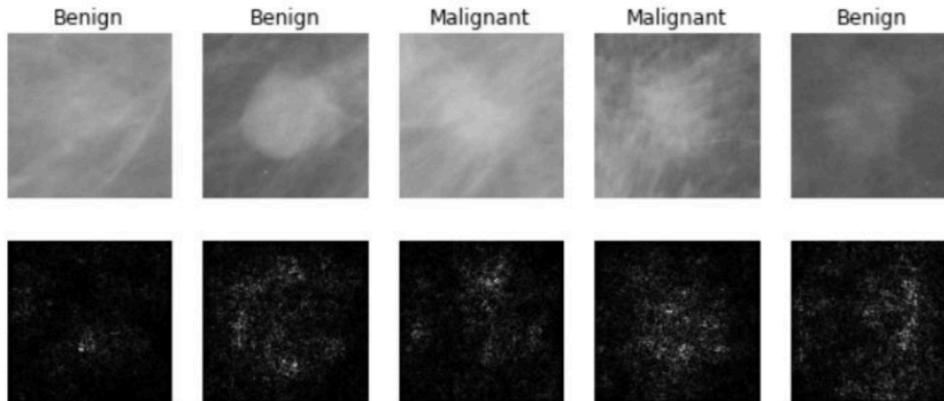


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

“AlexNet”

CNN is everywhere

Yes ...



[Levy et al. 2016]

Figure copyright Levy et al. 2016.
Reproduced with permission.



[Dieleman et al. 2014]

From left to right: [public domain by NASA](#), usage permitted by
ESA/Hubble, [public domain by NASA](#), and [public domain](#).



[Sermanet et al. 2011]

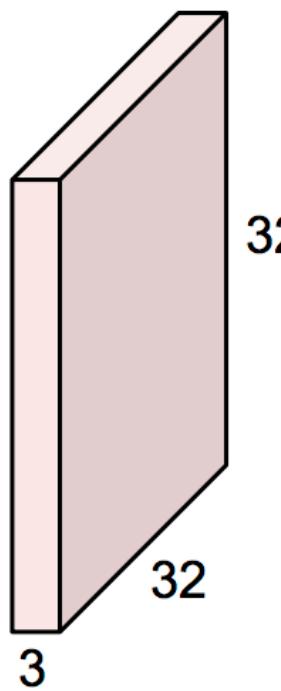
[Ciresan et al.]

Photos by Lane McIntosh.
Copyright CS231n 2017.

CNN simple

One simple step of convolution:

32x32x3 image



5x5x3

