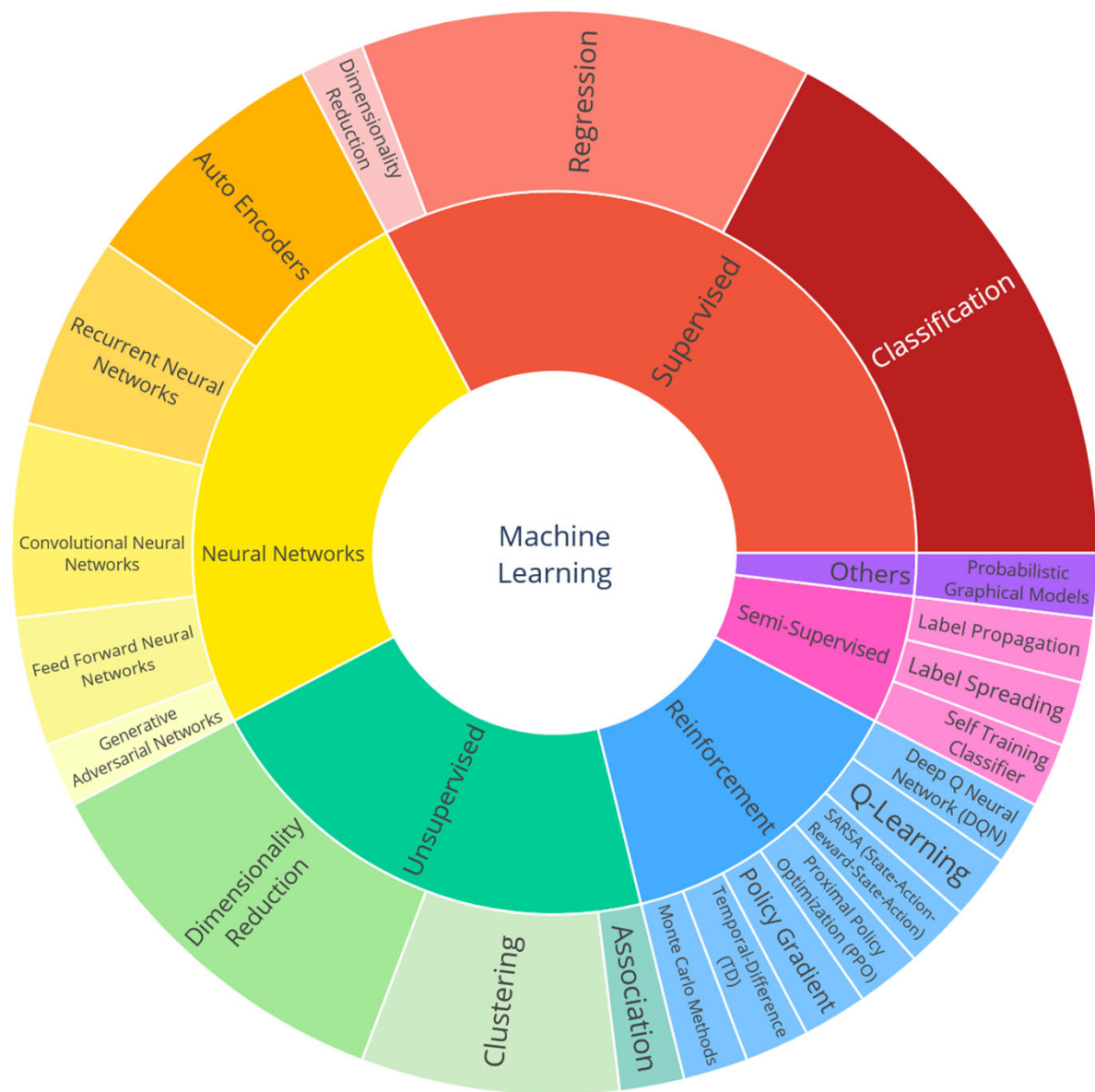




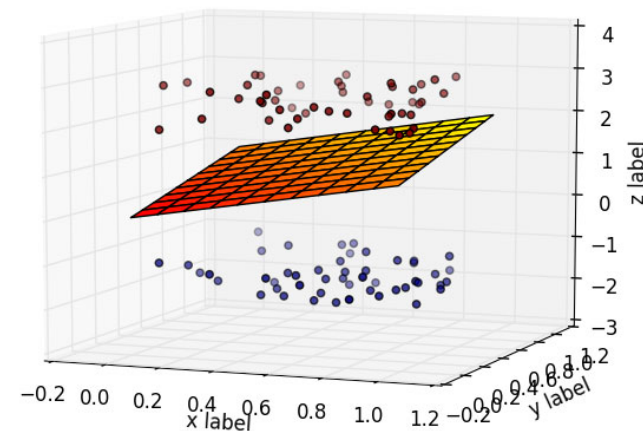
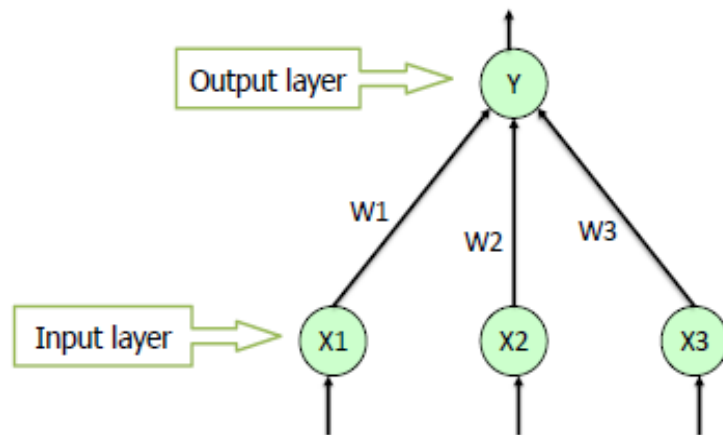
Introduction to Machine Learning

Neural Networks

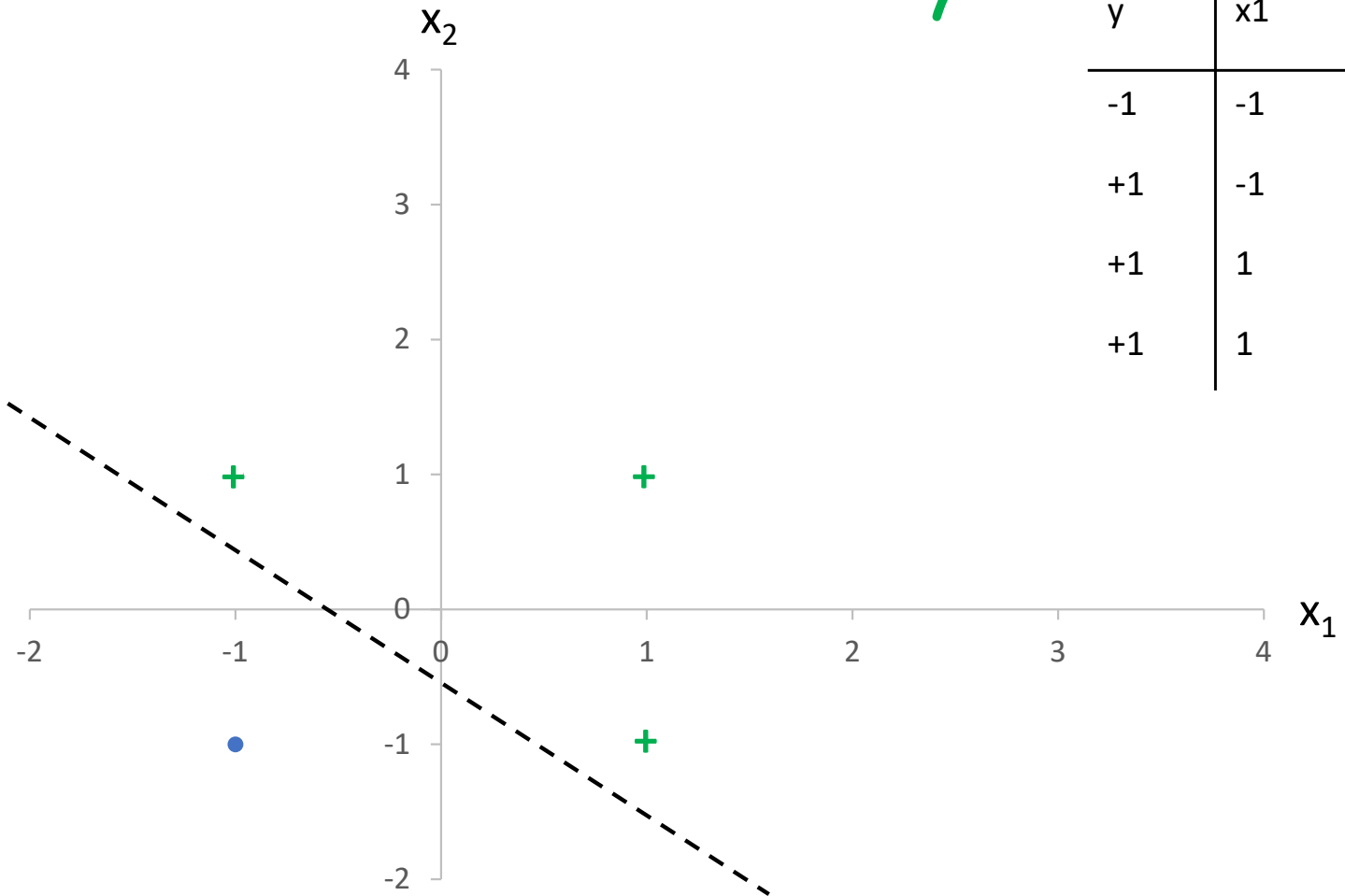


Linearly separable

- If the classes can be separated by a hyperplane, then they are linearly separable

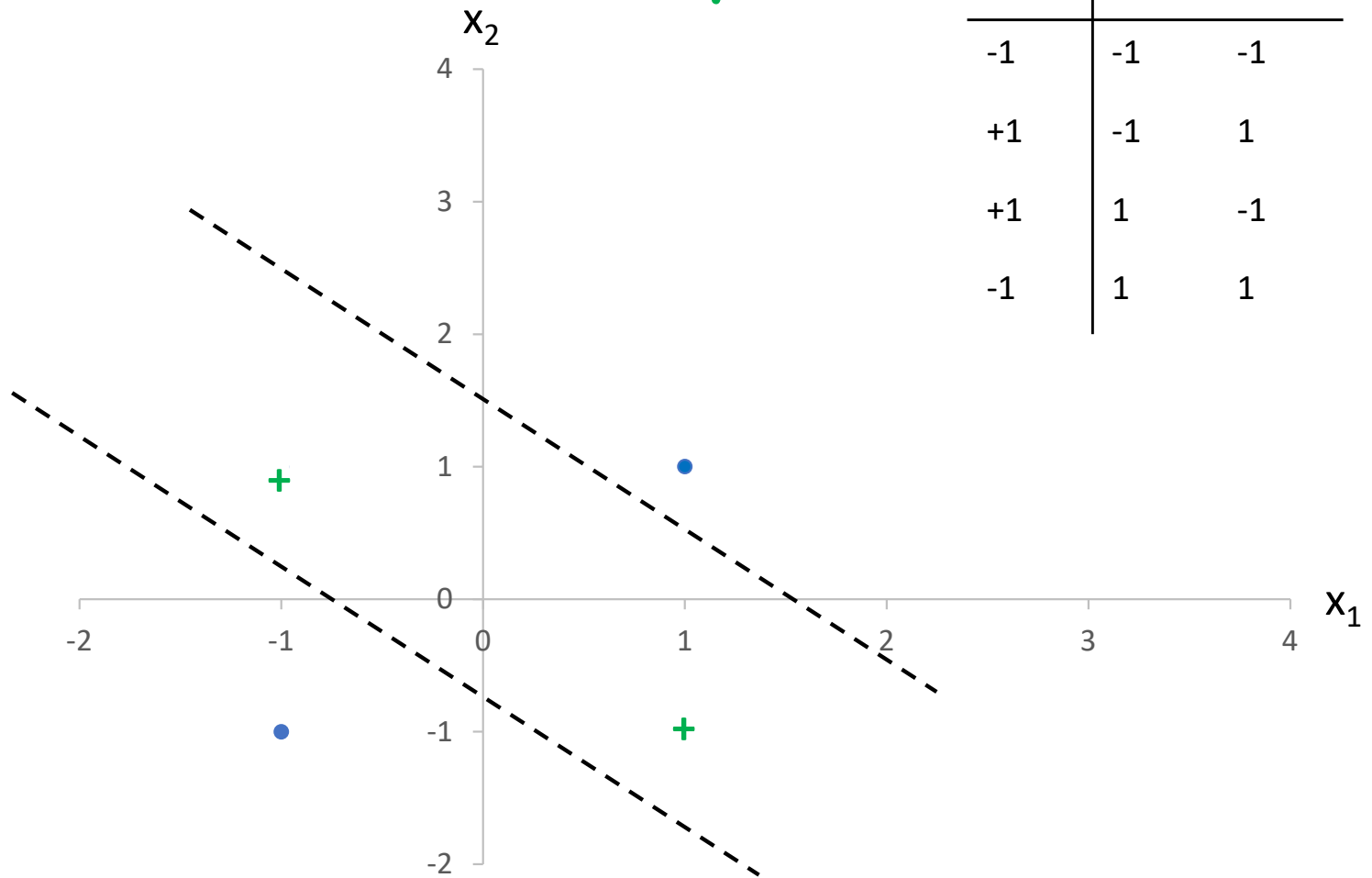


Logical OR decision boundary

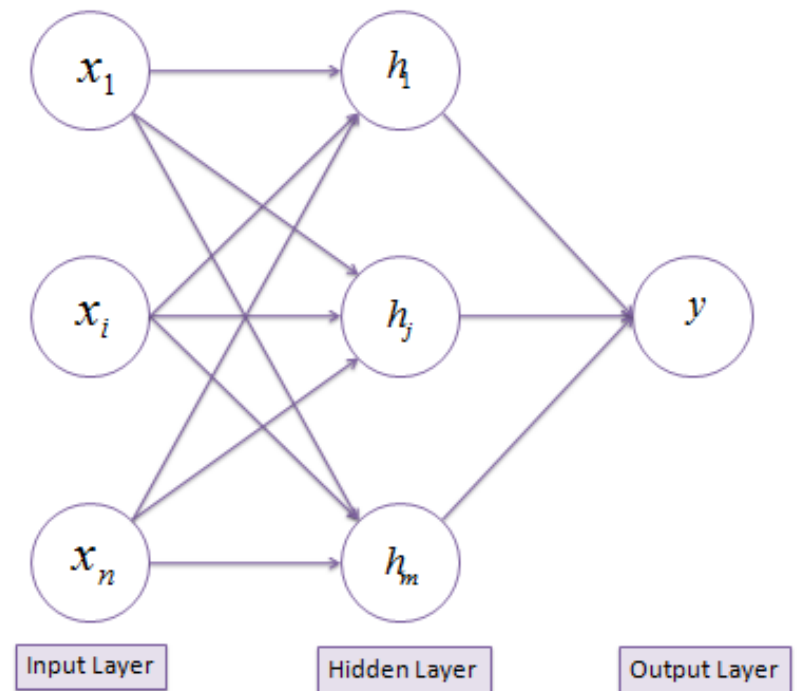
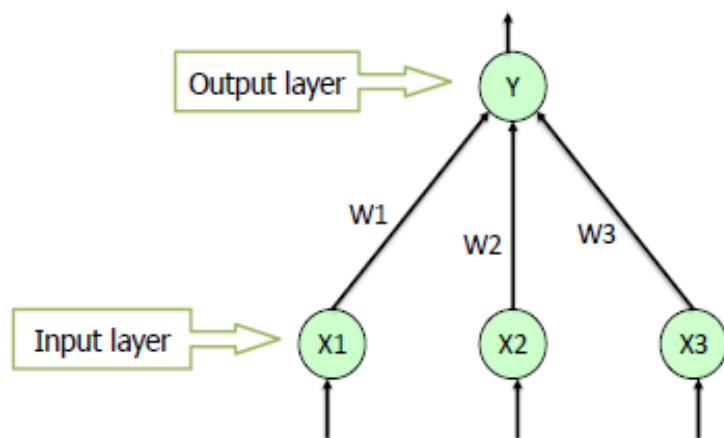


XOR decision boundary

y	x_1	x_2
-1	-1	-1
+1	-1	1
+1	1	-1
-1	1	1

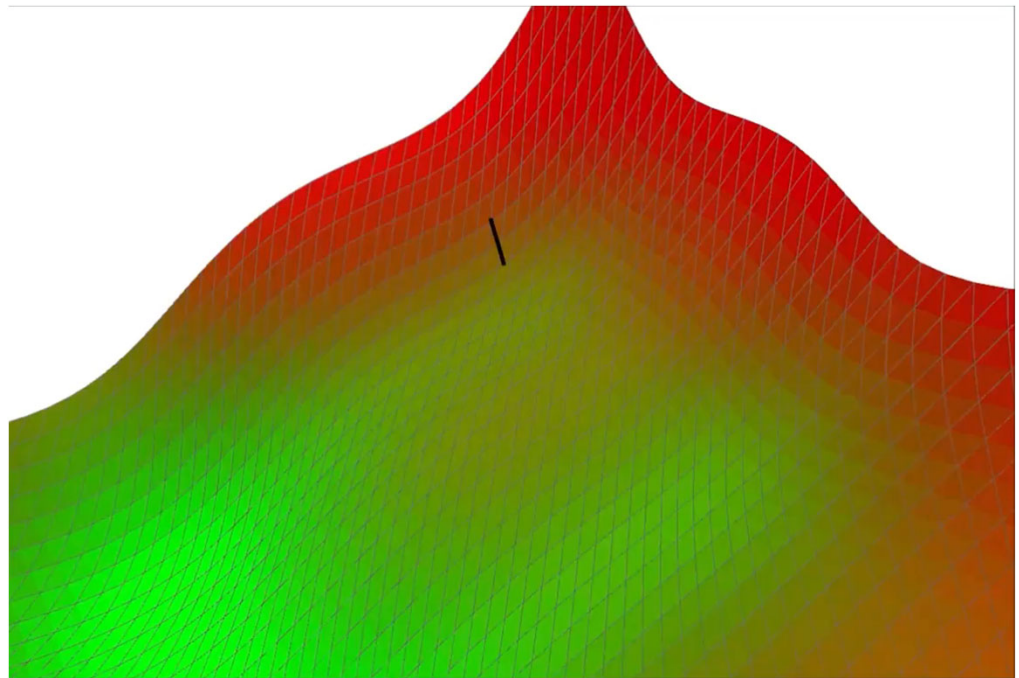


Adding layers



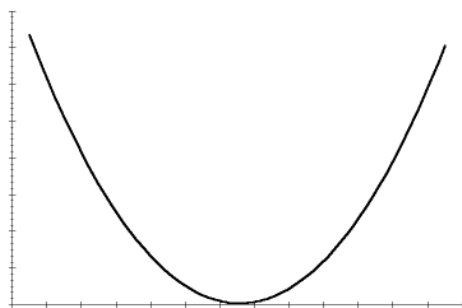
Activation function and gradient descent

$$h_i = f(w_i \cdot x)$$

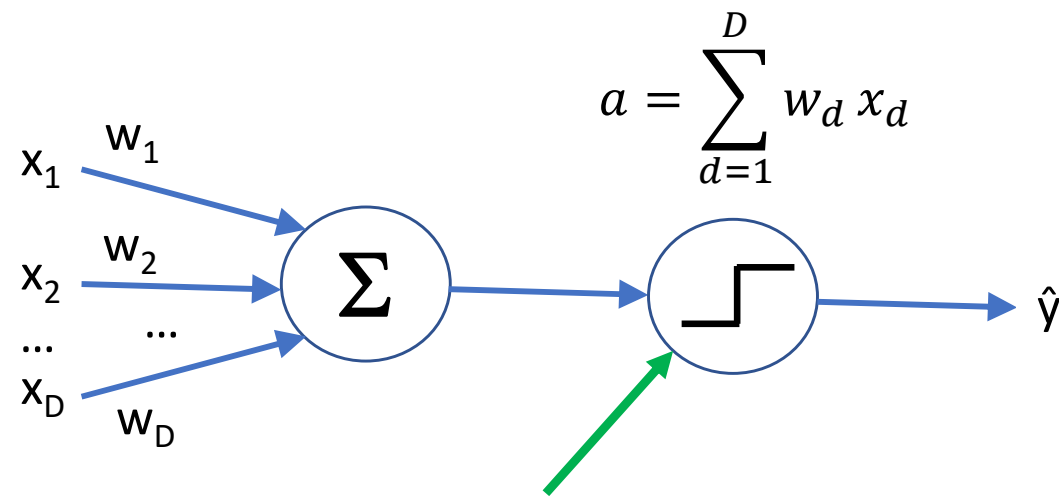


Algorithm 21 GRADIENTDESCENT($\mathcal{F}, K, \eta_1, \dots$)

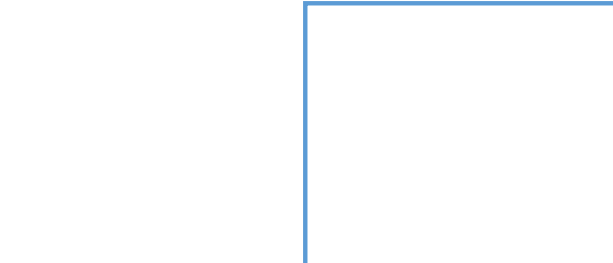
```
1:  $\mathbf{z}^{(0)} \leftarrow \langle 0, 0, \dots, 0 \rangle$  // initialize variable we are optimizing
2: for  $k = 1 \dots K$  do
3:    $\mathbf{g}^{(k)} \leftarrow \nabla_{\mathbf{z}} \mathcal{F}|_{\mathbf{z}^{(k-1)}}$  // compute gradient at current location
4:    $\mathbf{z}^{(k)} \leftarrow \mathbf{z}^{(k-1)} - \eta^{(k)} \mathbf{g}^{(k)}$  // take a step down the gradient
5: end for
6: return  $\mathbf{z}^{(K)}$ 
```



Perceptron neuron

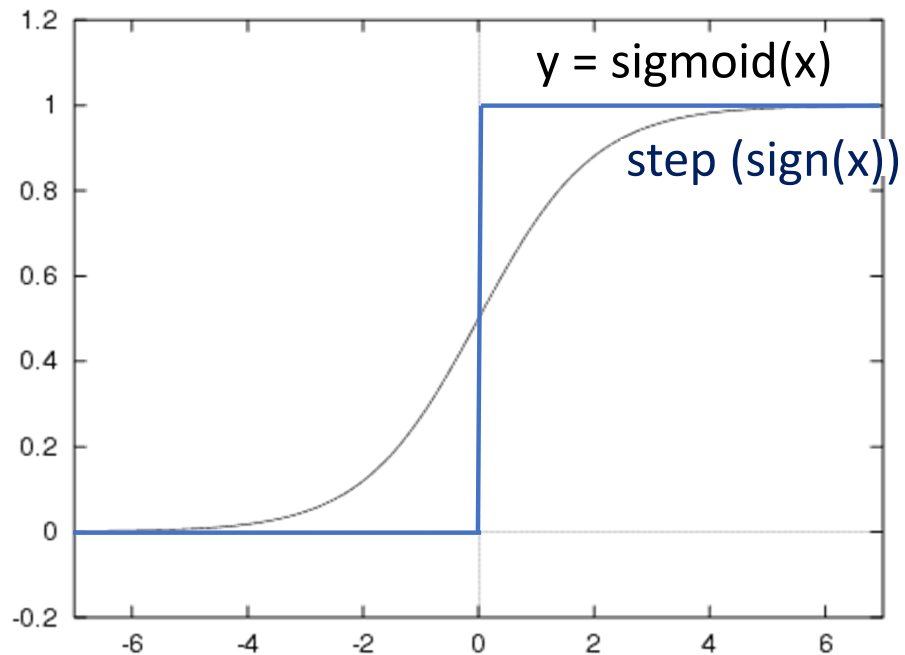


Not differentiable!

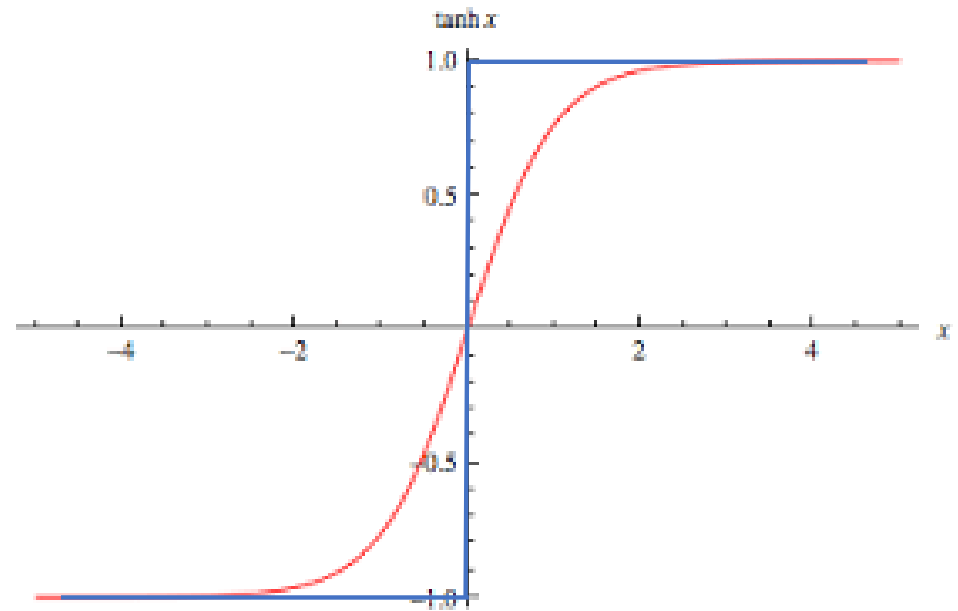


derivative = infinity here

Alternatives that are differentiable



We will use this one, $\tanh(x)$



Derivative is $1 - \tanh^2(x)$

Algorithm 25 `TWO_LAYER_NETWORK_PREDICT`(W, v, \hat{x})

```
1: for  $i = 1$  to number of hidden units do  
2:    $h_i \leftarrow \tanh(w_i \cdot \hat{x})$  // compute activation of hidden unit  $i$   
3: end for  
4: return  $v \cdot h$  // compute output unit
```

- W is matrix of weights from input nodes to hidden nodes
- v is vector of weights from hidden nodes to output node
- Note the different way the hidden value is computed in comparison with the output value

XOR

y	x1	x2
-1	-1	-1
+1	-1	1
+1	1	-1
-1	1	1

Top hidden node computes “or”, weights=1.0, bias=0.5.

Associate bias with each node (handled various ways, here listed in circle).

$$x_1=-1, x_2=-1, \tanh(1*-1 + 1*-1 - 0.5) = -0.905 \quad (-1)$$

$$x_1=-1, x_2=1, \tanh(1*-1 + 1*1 - 0.5) = +0.462 \quad (1)$$

$$x_1=1, x_2=-1, \tanh(1*1 + 1*-1 - 0.5) = +0.462 \quad (1)$$

$$x_1=1, x_2=1, \tanh(1*1 + 1*1 - 0.5) = +0.987 \quad (1)$$

Bottom hidden node computes “nand”, weights=-1.0, bias=1.5.

$$x_1=-1, x_2=-1, \tanh(-1*-1 + -1*-1 - 1.5) = +0.998 \quad (1)$$

$$x_1=-1, x_2=1, \tanh(-1*-1 + -1*1 - 1.5) = +0.905 \quad (1)$$

$$x_1=1, x_2=-1, \tanh(-1*1 + -1*-1 - 1.5) = +0.905 \quad (1)$$

$$x_1=1, x_2=1, \tanh(-1*1 + -1*1 - 1.5) = -0.462 \quad (-1)$$

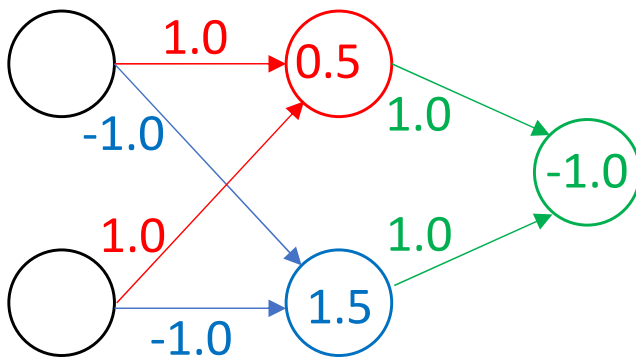
Output computes “and(or,nand)”, weights=1.0, bias=1.0.

$$x_1=-1, x_2=-1, h = (-.905, .998) = -.907 \quad (-1)$$

$$x_1=-1, x_2=1, h = (.462, .905) = .367 \quad (1)$$

$$x_1=1, x_2=-1, h = (.462, .905) = .367 \quad (1)$$

$$x_1=1, x_2=1, h = (.987, -.462) = -.475 \quad (-1)$$



How big of a network do I need?

Let F be a continuous function on a bounded subset of D -dimensional space. Then there exists a two-layer neural network \hat{F} with a finite number of hidden units that approximate F arbitrarily well. Namely, for all \mathbf{x} in the domain of F , $|F(\mathbf{x}) - \hat{F}(\mathbf{x})| < \varepsilon$.

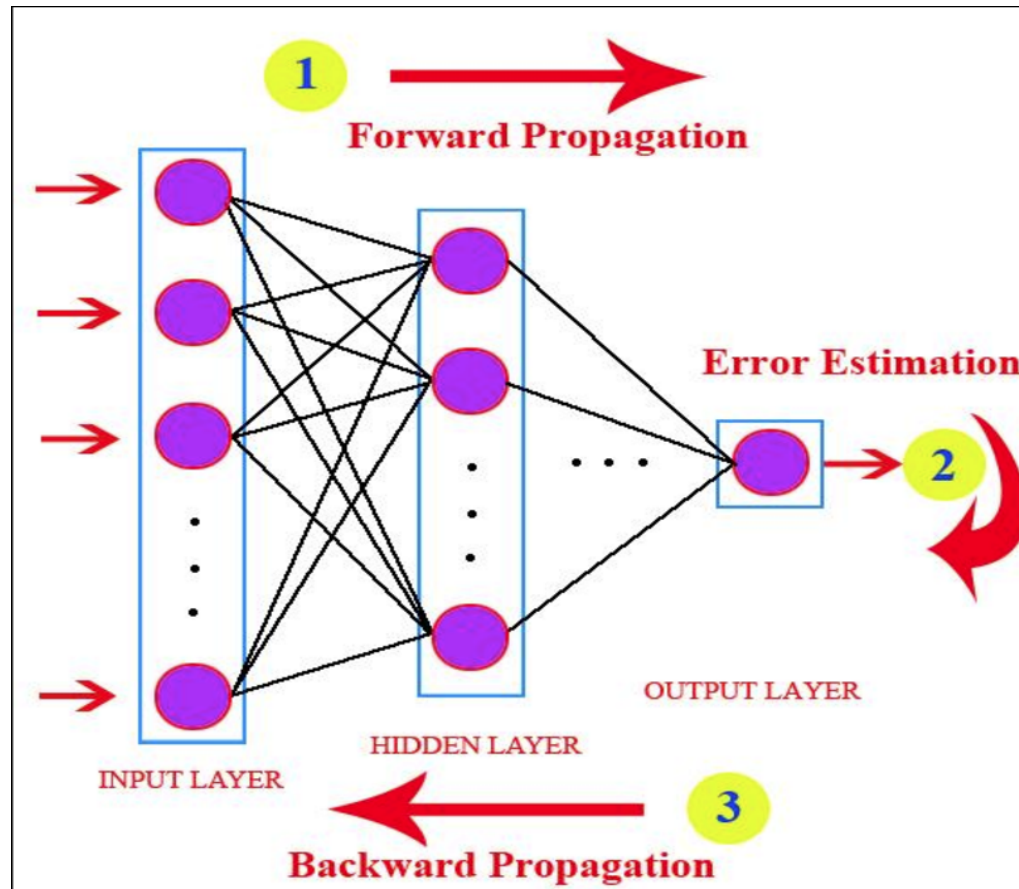
How many hidden units?

If D dimensions, K hidden units, then $(D+2)K$ parameters

1 parameter for bias, 1 for weight to the output node, thus $N = (D+2)K$, $N/(D+2)=K$

If you want 1-2 examples for each parameter you are learning, then use

$$\text{\#hidden nodes} = K = \left\lceil \frac{N}{D} \right\rceil$$



Backpropagation

- We know how to compute output
- How do we learn weights?

backpropagation = gradient descent + chain rule

Review: Linear Regression

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) &= \frac{\partial}{\partial \theta_j} \bullet \frac{1}{2n} \sum_{i=1}^n (h(x^i) - y^i)^2 \\ &= \frac{\partial}{\partial \theta_j} \bullet \frac{1}{2n} \sum_{i=1}^n (\theta_0 + \theta_1 x^i - y^i)^2 \end{aligned}$$

Repeat until convergence {

$$\theta_j := \theta_j + \alpha \sum_{i=1}^n (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)} \quad (\text{for every } j)$$

}

Optimization function

- Based on loss function: Squared error

$$\min_{W, v} \sum_n \frac{1}{2} \left(y_n - \sum_i v_i f(w_i \cdot x_n) \right)^2$$

Adjusting weights to output node

- Differentiate objective with respect to v

$$\min_{\mathbf{W}, v} \sum_n \frac{1}{2} \left(y_n - \sum_i v_i f(\mathbf{w}_i \cdot \mathbf{x}_n) \right)^2$$

$$\nabla_v = - \sum_n e_n \mathbf{h}_n$$

Adjusting weights to hidden node

- Only compensate for *portion* of error for that hidden node
- Not trying to produce specific output value

Adjusting weights to hidden node

$$\mathcal{L}(\mathbf{W}) = \frac{1}{2} \left(y - \sum_i v_i f(w_i \cdot x) \right)^2$$

$$\frac{\partial \mathcal{L}}{\partial w_i} = \frac{\partial \mathcal{L}}{\partial f_i} \frac{\partial f_i}{\partial w_i}$$

$$\frac{\partial \mathcal{L}}{\partial f_i} = - \left(y - \sum_i v_i f(w_i \cdot x) \right) v_i = -ev_i$$

$$\frac{\partial f_i}{\partial w_i} = f'(w_i \cdot x)x$$

This is the **gradient** with respect to w_i

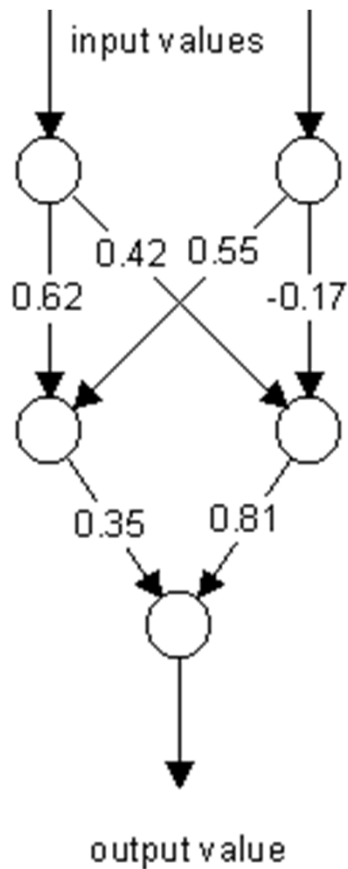
- (network error) * (weight from this hidden node to output node) *
(derivative of activation function) * (feature value for this data point x)

$$\begin{aligned} \nabla_{w_i} &= -ev_i f'(w_i \cdot x)x \\ &= -ev_i (1 - \tanh^2(a_i))x \end{aligned}$$

Algorithm 26 `TWO_LAYER_NETWORK_TRAIN`($\mathbf{D}, \eta, K, \text{MaxIter}$)

```
1:  $\mathbf{W} \leftarrow D \times K$  matrix of small random values // initialize input layer weights
2:  $\mathbf{v} \leftarrow K$ -vector of small random values // initialize output layer weights
3: for  $iter = 1 \dots \text{MaxIter}$  do
4:    $\mathbf{G} \leftarrow D \times K$  matrix of zeros // initialize input layer gradient
5:    $\mathbf{g} \leftarrow K$ -vector of zeros // initialize output layer gradient
6:   for all  $(x, y) \in \mathbf{D}$  do
7:     for  $i = 1$  to  $K$  do
8:        $a_i \leftarrow \mathbf{w}_i \cdot \hat{\mathbf{x}}$ 
9:        $h_i \leftarrow \tanh(a_i)$  // compute activation of hidden unit  $i$ 
10:    end for
11:     $\hat{y} \leftarrow \mathbf{v} \cdot \mathbf{h}$  // compute output unit
12:     $e \leftarrow y - \hat{y}$  // compute error
13:     $\mathbf{g} \leftarrow \mathbf{g} - e\mathbf{h}$  // update gradient for output layer
14:    for  $i = 1$  to  $K$  do
15:       $\mathbf{G}_i \leftarrow \mathbf{G}_i - ev_i(1 - \tanh^2(a_i))x$  // update gradient for input layer
16:    end for
17:  end for
18:   $\mathbf{W} \leftarrow \mathbf{W} - \eta \mathbf{G}$  // update input layer weights
19:   $\mathbf{v} \leftarrow \mathbf{v} - \eta \mathbf{g}$  // update output layer weights
20: end for
21: return  $\mathbf{W}, \mathbf{v}$ 
```

Example



Example Input \rightarrow Target: 0 1 \rightarrow 0; Assume learning rate=1.0, ignore bias

Forward propagation:

Dot product for left hidden node = $0 \cdot 0.62 + 1 \cdot 0.55 = 0.55$

Dot product for right hidden node = $0 \cdot 0.42 + 1 \cdot -0.17 = -0.17$

$h(\text{left}) = \tanh(.55) = .5005$; $h(\text{right}) = \tanh(-.17) = -.1684$

$\text{output} = 0.35 \cdot .5005 + 0.81 \cdot -.1684 = 0.039$; $\text{error} = 0 - 0.039 = -0.039$

Backward propagation:

$g = \text{error} \cdot h$, $g(\text{left}) = -0.039 \cdot .5005 = .0195$; $v(\text{left}) = 0.35 \cdot .0195 = 0.33$

$g(\text{right}) = -0.039 \cdot -.1684 = -.007$; $v(\text{right}) = 0.81 \cdot -.007 = 0.82$

$W(\text{left_to_left}) += -.039 \cdot .35 \cdot (1 - \tanh^2(.55))0 = .62 + 0$ ($W = 0.62$, the same)

$W(\text{left_to_right}) += -.039 \cdot .81 \cdot (1 - \tanh^2(-.17))0 = .42 + 0$ ($W = 0.42$, the same)

$W(\text{right_to_left}) += -.039 \cdot .35 \cdot (1 - \tanh^2(.55))1 = 0.55 + (0.39 \cdot .35 \cdot .75 \cdot 1) = .52$;

$W(\text{right_to_right}) += -.039 \cdot .81 \cdot (1 - \tanh^2(-.17))1 = -.17 + (0.39 \cdot .81 \cdot .97 \cdot 1) = -.20$

Let's try this out

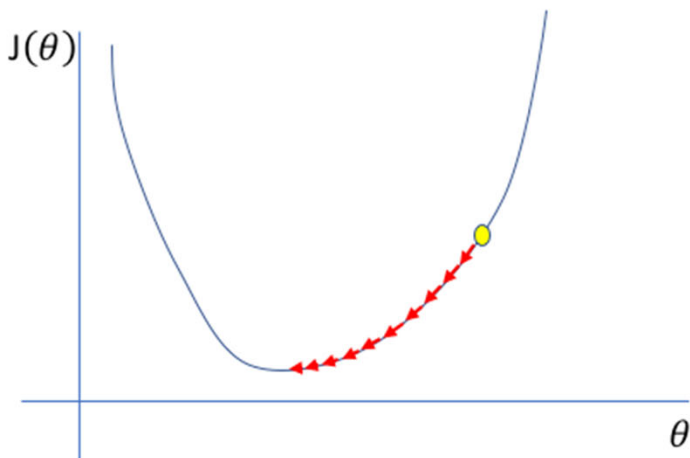
Practical issues

- Initializing weights
- Why not initialize to 0? Can get stuck in local optimum.
- If $W=0$ and $v=0$ then the activation h_i of hidden units will all be 0 (because $W=0$). Because h will be 0 $e \cdot h$ will be 0 so the weights v will not change.
- Note if sigmoid is used instead of tanh then activation will be non-zero so weights will change but they will all change identically and hidden unit values will always be the same. The model will eventually converge but hidden nodes are ignored.
- Moral: neural networks are sensitive to their initialization. Random (small) initialization is most effective.

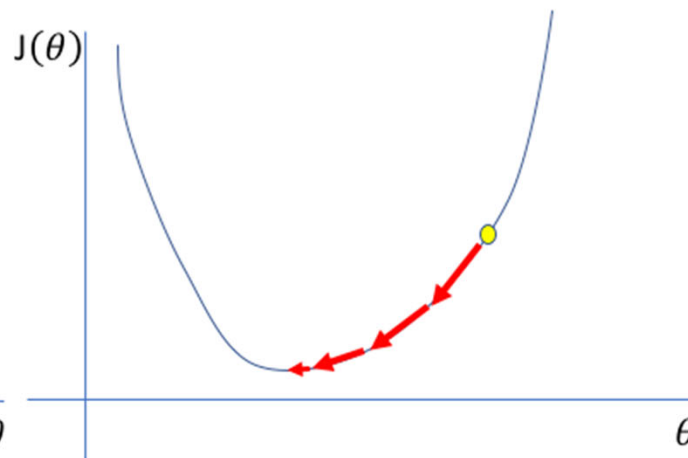
```
1:  $W \leftarrow D \times K$  matrix of small random values
2:  $v \leftarrow K$ -vector of small random values
3: for  $iter = 1 \dots MaxIter$  do
4:    $G \leftarrow D \times K$  matrix of zeros
5:    $g \leftarrow K$ -vector of zeros
6:   for all  $(x, y) \in D$  do
7:     for  $i = 1$  to  $K$  do
8:        $a_i \leftarrow w_i \cdot \hat{x}$ 
9:        $h_i \leftarrow \tanh(a_i)$ 
10:    end for
11:     $\hat{y} \leftarrow v \cdot h$ 
12:     $e \leftarrow y - \hat{y}$ 
13:     $g \leftarrow g - eh$ 
14:    for  $i = 1$  to  $K$  do
15:       $G_i \leftarrow G_i - ev_i(1 - \tanh^2(a_i))x$ 
16:    end for
17:  end for
18:   $W \leftarrow W - \eta G$ 
19:   $v \leftarrow v - \eta g$ 
20: end for
21: return  $W, v$ 
```

Learning rate

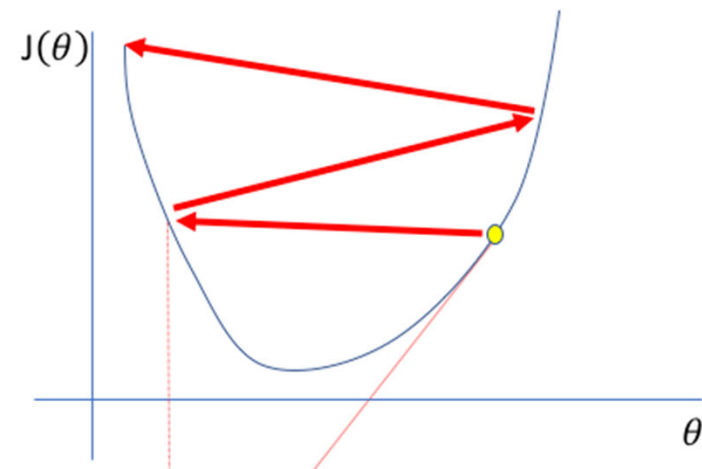
Too low



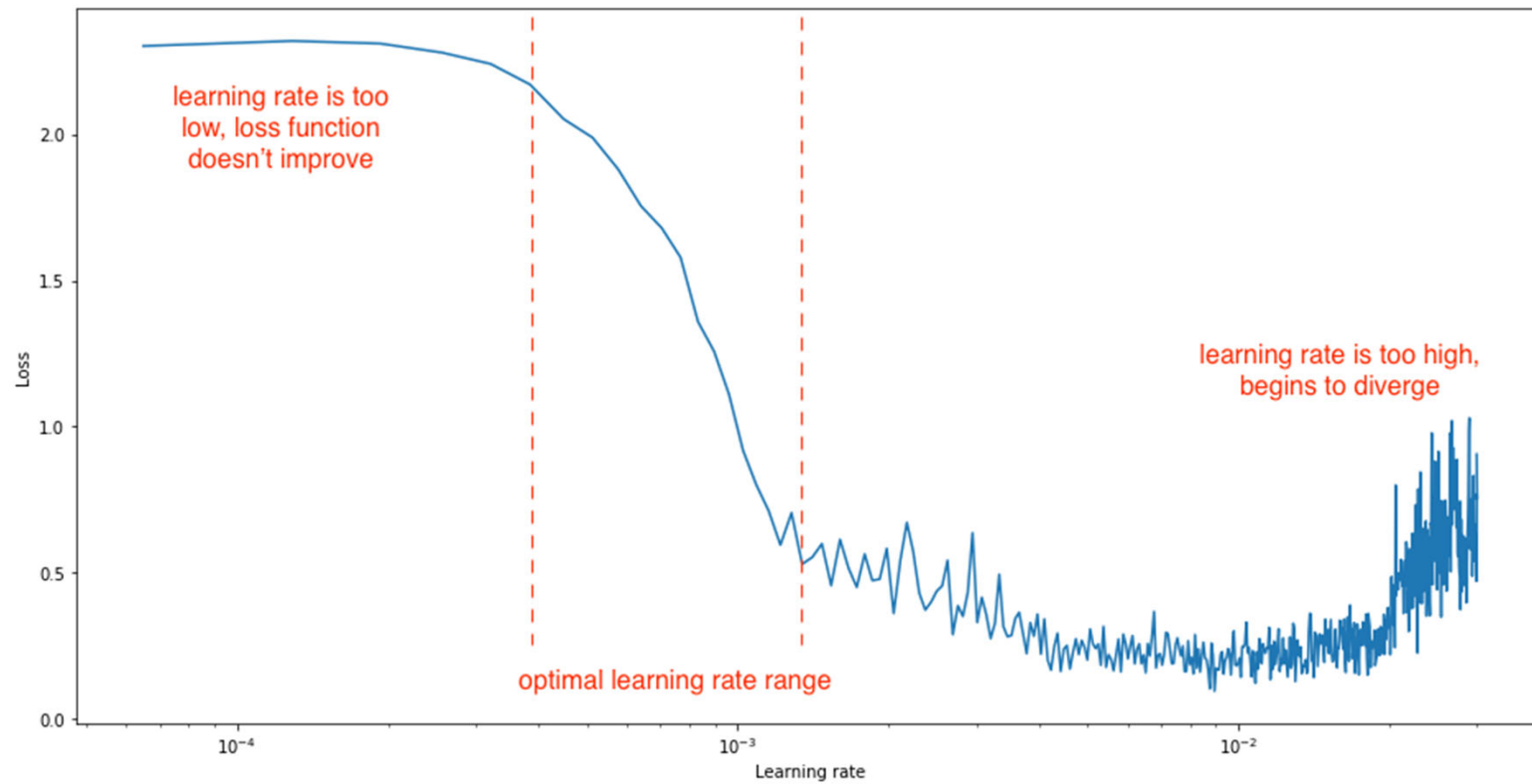
Just right



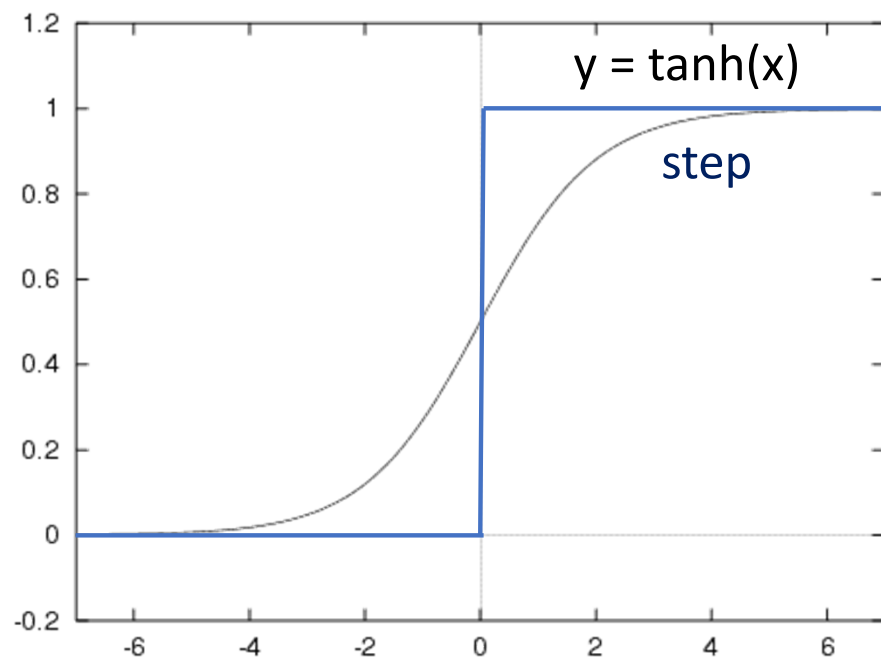
Too high




Learning rate

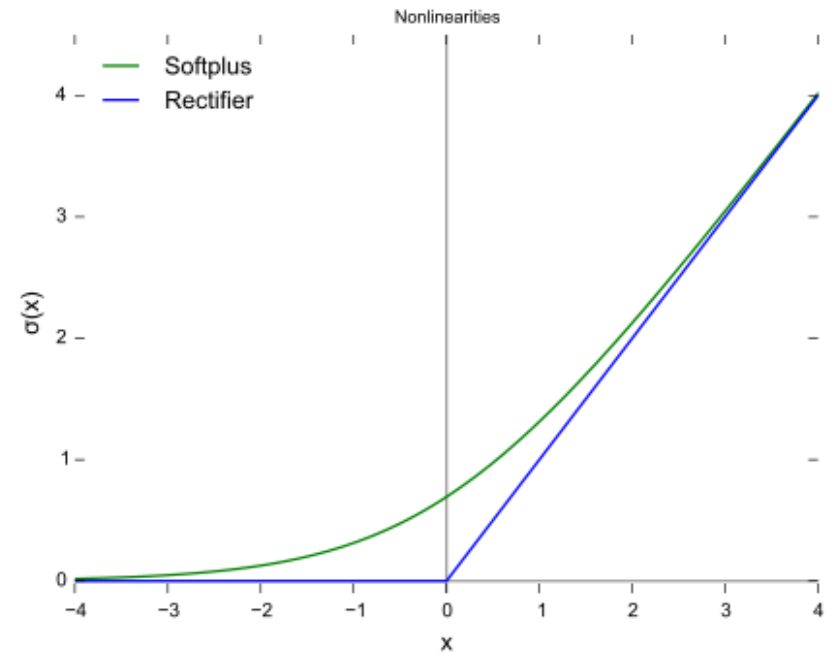
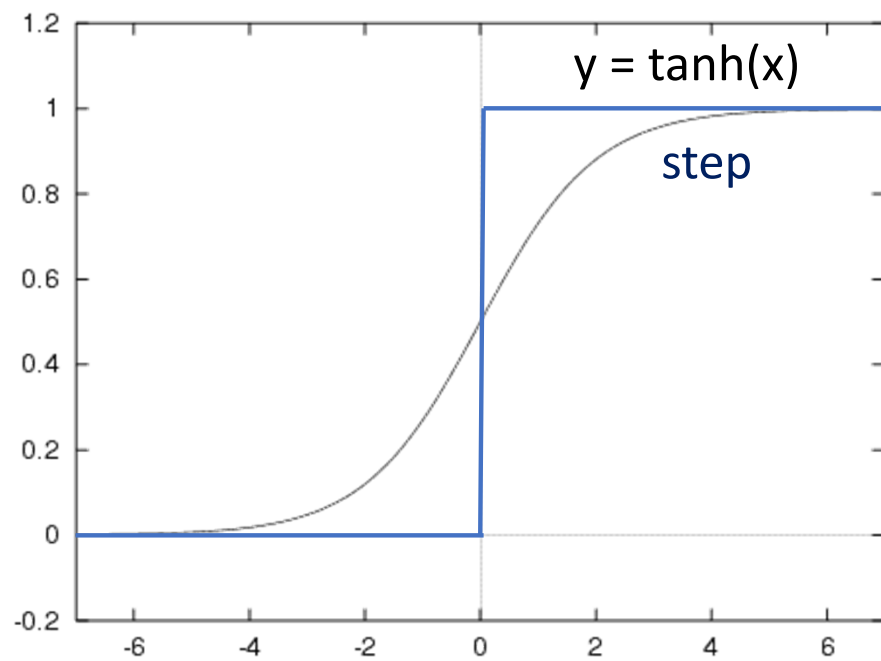


Activation function



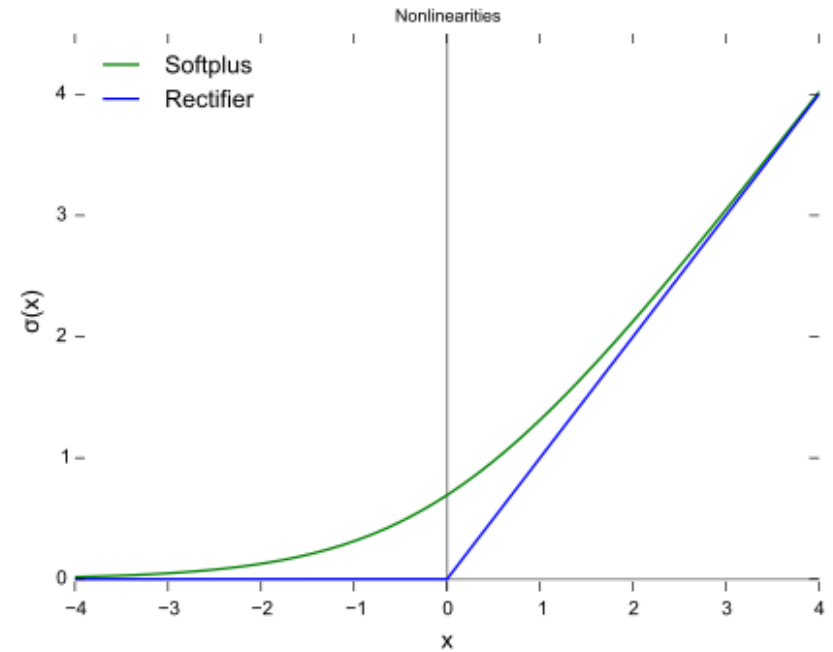
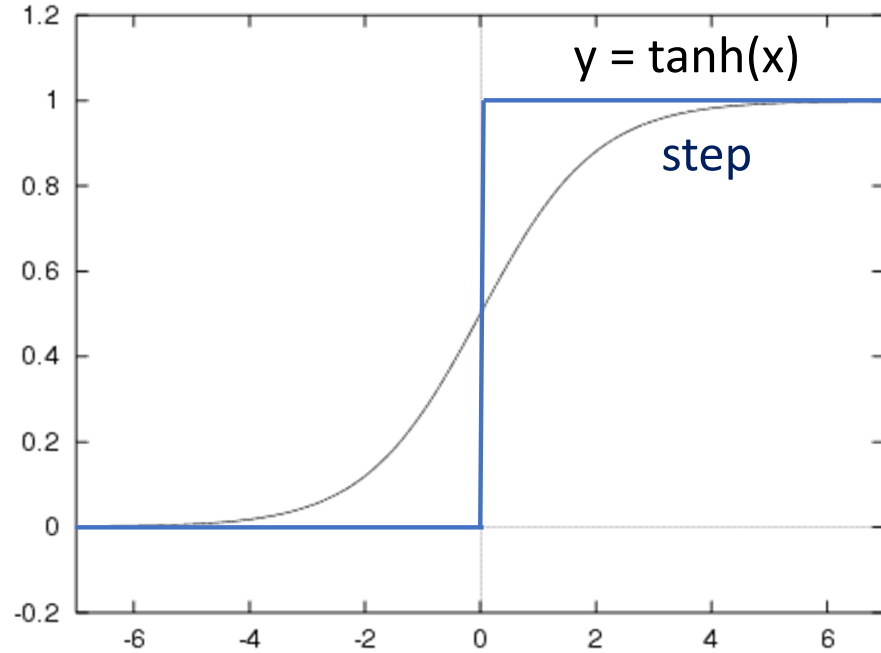
ReLU

	$f(x) = \begin{cases} 0 & \text{for } x \leq 0 \\ x & \text{for } x > 0 \end{cases} = \max\{0, x\} = x \mathbf{1}_{x>0}$	$f'(x) = \begin{cases} 0 & \text{for } x \leq 0 \\ 1 & \text{for } x > 0 \end{cases}$	$[0, \infty)$
---	--	---	---------------

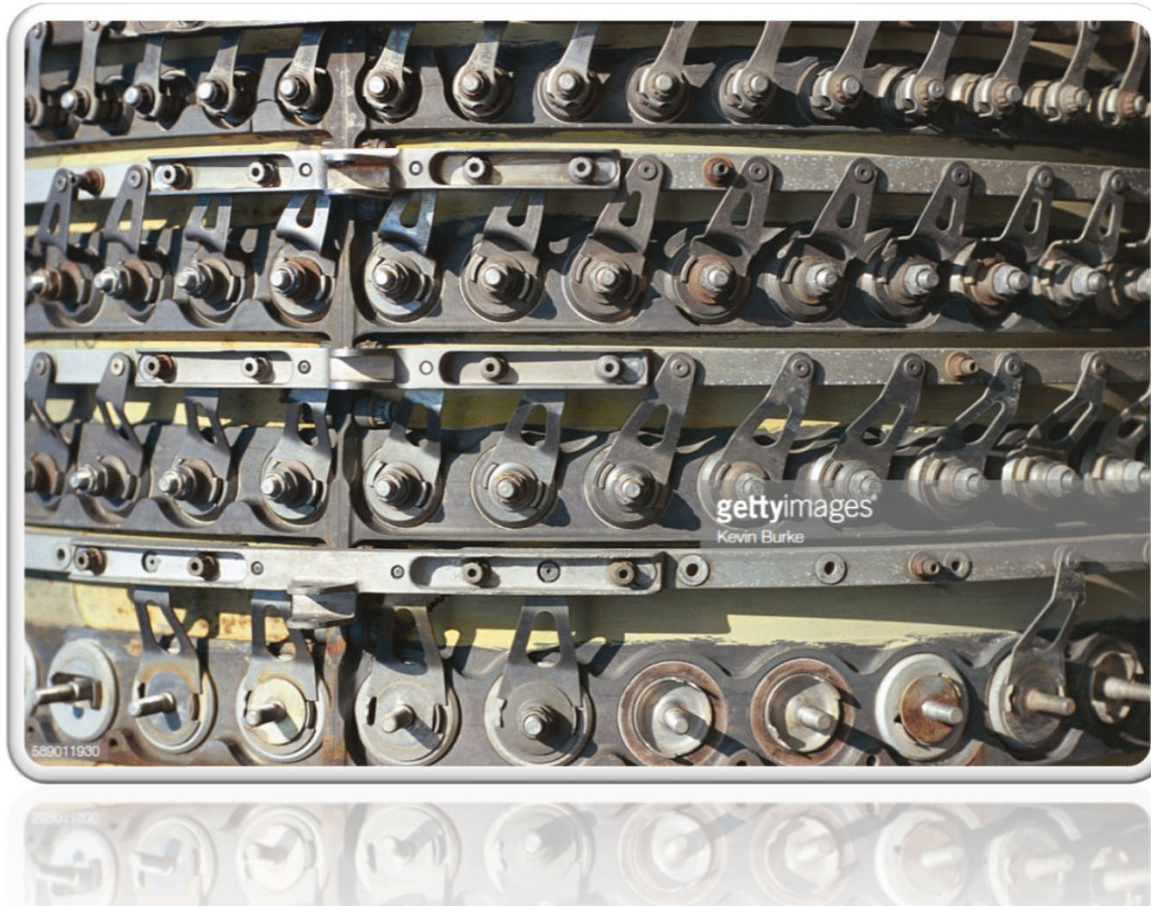


Softplus

$$f(x) = \log(1 + e^x)$$



Hyperparameters



Practical issues

- Multiple classes
- Softmax

Class	Probability
apple	0.001
bear	0.040
candy	0.008
dog	0.950
egg	0.001

