



THE STATE UNIVERSITY
OF NEW JERSEY

Machine Learning

Basic Algorithms

Edgar Granados

R

Classification

- Set of training examples: $\underline{X_{train}} = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, \underline{y_n})\}$
- Learn: $f(x_i) = y_i$
- Evaluate performance: $X_{test} = \{(x_{n+1}, y_{n+1}), \dots, (x_m, y_m)\}$
 \hookrightarrow Net knows while training

R

Naive Bayes Algorithm

- Recall Bayes Rule:

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

- Likelihood:

$$\cdot L(x) = \frac{P(y = \text{true} | x)}{P(y = \text{false} | x)} = \frac{P(x | y = \text{true}) P(y = \text{true})}{P(x | y = \text{false}) P(y = \text{false})}$$

if $L(x) > 1$
 y is true

if $L(x) < 1$
 y is false

R

Naive Bayes Algorithm

- Estimate $P(y)$ from the training examples

$$\bullet P(\text{y} = \cancel{\text{true}}) = \frac{\# \text{times } y_i = \text{true in } X_{\text{train}}}{n}$$

$$\bullet P(\text{y} = \cancel{\text{false}}) = \frac{\# \text{times } y_i = \text{false in } X_{\text{train}}}{n}$$

R

Naive Bayes Algorithm

x_{train}

- From the training examples
1. Define ϕ , a function that maps each input x into a features vector

$$\phi(x) = (\phi_1(x), \phi_2(x), \dots, \phi_f(x))$$

1

R

Naive Bayes Algorithm

- From the training examples
 1. Define ϕ , a function that maps each input x into a features vector
 2. Estimate the probability distribution of each feature ϕ_j

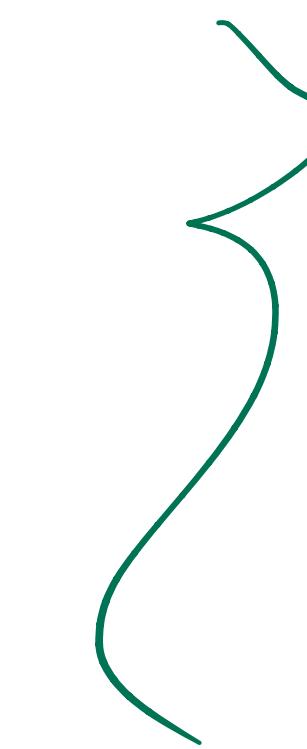
$$P(\phi_j(x) = v | y = \text{true}) = \frac{\# \phi_j(x) = v \text{ AND } y = \text{true} \in X_{\text{train}}}{\# y = \text{true} \in X_{\text{train}}}$$

$$P(\phi_j(x) = v | y = \text{false}) = \frac{\# \phi_j(x) = v \text{ AND } y = \text{false} \in X_{\text{train}}}{\# y = \text{false} \in X_{\text{train}}}$$

R

Naive Bayes Algorithm

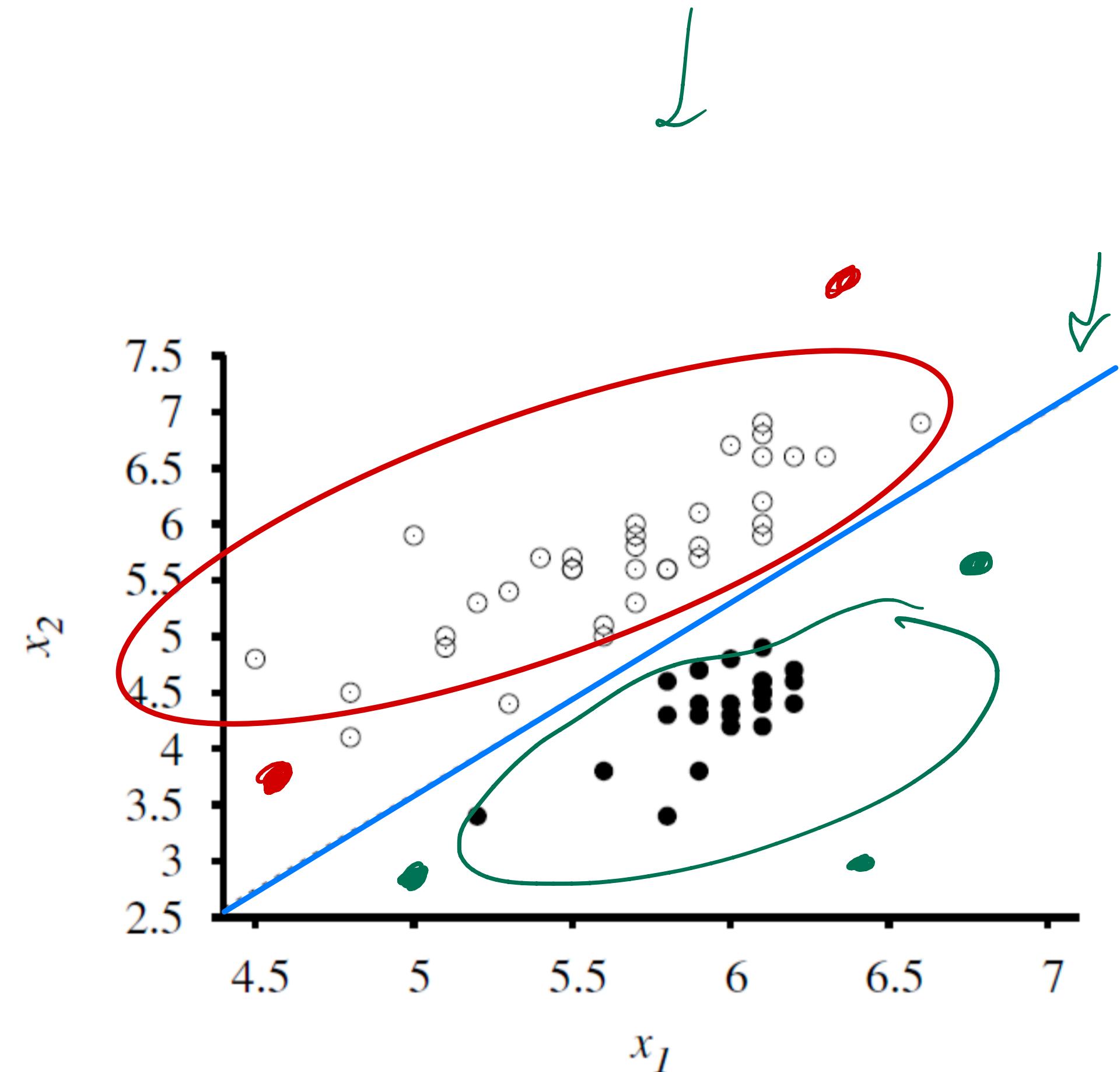
- From the training examples
 1. Define ϕ , a function that maps each input x into a features vector
 2. Estimate the probability distribution of each feature ϕ_j
 3. Estimate $P(x | y)$ using the naive Bayes assumption


$$p(x|y = \text{true}) = \prod_{j=1}^l p(\phi_j(x)|y = \text{true})$$
$$p(x|y = \text{false}) = \prod_{j=1}^l p(\phi_j(x)|y = \text{false})$$

R

Linear Classifier

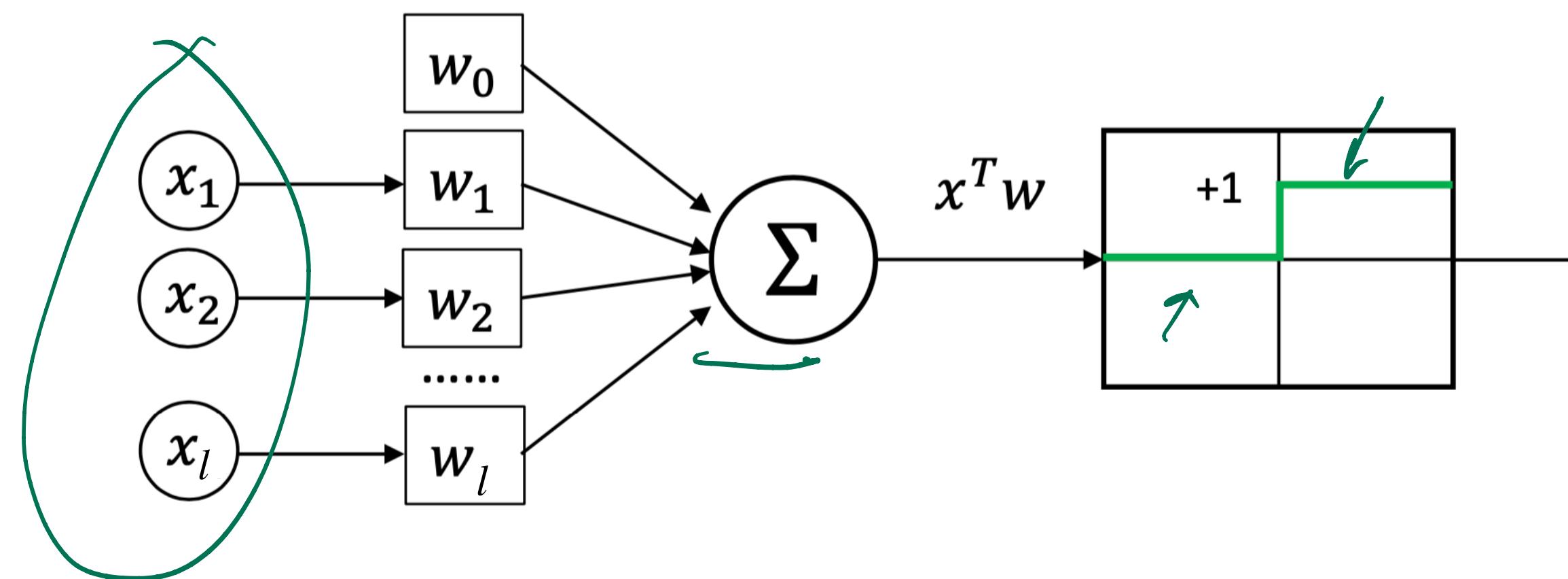
- A linear classifier for l features is an $(l - 1)$ -dimensional plane separating the positive/negative labels



R

Linear Classifier

Perceptron



- Learning:
- Find the appropriate weights

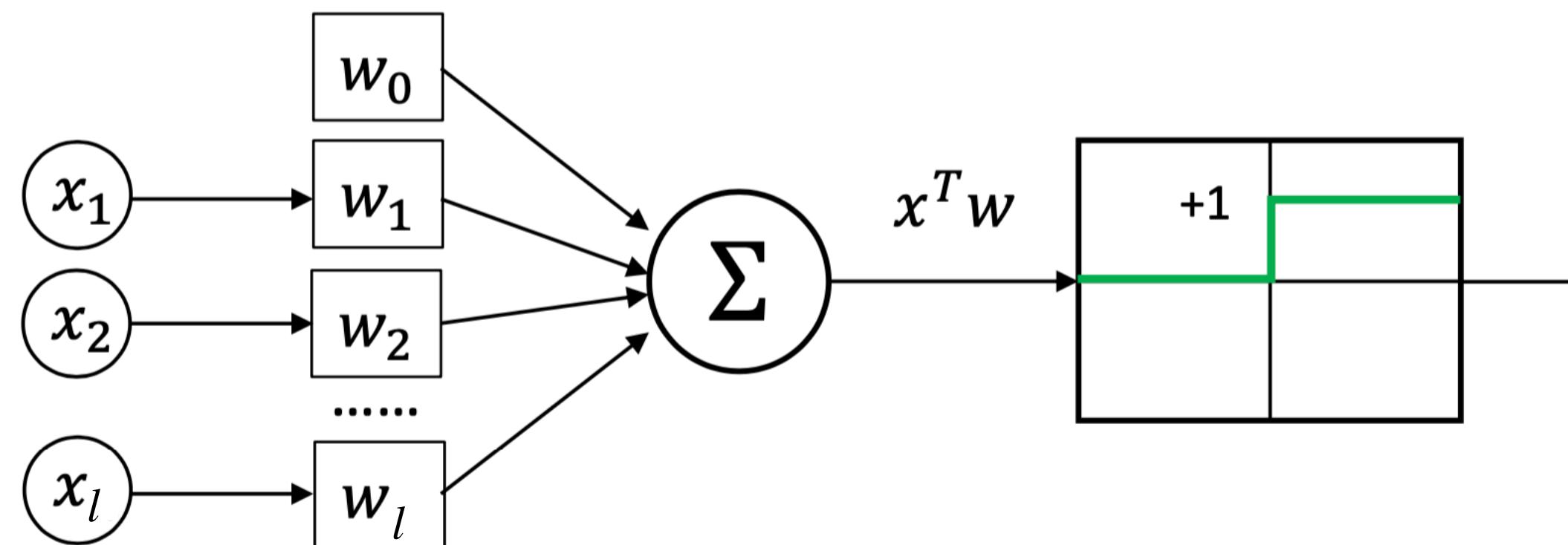
$$y = \begin{cases} 1, & w_0 + w_1\phi_1(x_1) + \dots + w_l\phi_l(x_l) \geq 0 \\ 0, & w_0 + w_1\phi_1(x_1) + \dots + w_l\phi_l(x_l) < 0 \end{cases}$$

$y = Ax + b$
 $y = w_i x + w_0$

R

Linear Classifier

Perceptron



- Invented in 1957
- A single-layer neural network.
- Deep neural nets are nothing but Perceptrons stacked on top of each other

$$y = \begin{cases} 1, & w_0 + w_1\phi_1(x_1) + \dots + w_l\phi_l(x_l) \geq 0 \\ 0, & w_0 + w_1\phi_1(x_1) + \dots + w_l\phi_l(x_l) < 0 \end{cases}$$

R

Perceptron Algorithm

1. Initialize weights

R

Perceptron Algorithm

1. Initialize weights

2. $\forall (x_i, y_i) \in \underline{\mathcal{X}_{train}}$:

- $f(x_i, w) = w_0 + w_1\phi_1(x_i) + \dots + w_l\phi_l(x_i)$
- If $(f(x_i, w) \geq 0 \text{ AND } y_i = true)$ OR $(f(x_i, w) < 0 \text{ AND } y_i = false)$
 - Nothing... move to next example

Dots pt.

R

Perceptron Algorithm

1. Initialize weights
2. $\forall (x_i, y_i) \in \mathcal{X}_{train}:$
 - $f(x_i, w) = w_0 + w_1\phi_1(x_1) + \dots + w_n\phi_n(x_n)$
 - If $(f(x_i, w) \geq 0 \text{ AND } y_i = \text{true}) \text{ OR } (f(x_i, w) < 0 \text{ AND } y_i = \text{false}):$
 - Nothing... move to next example
 - Else: Update weights

R

Perceptron Algorithm

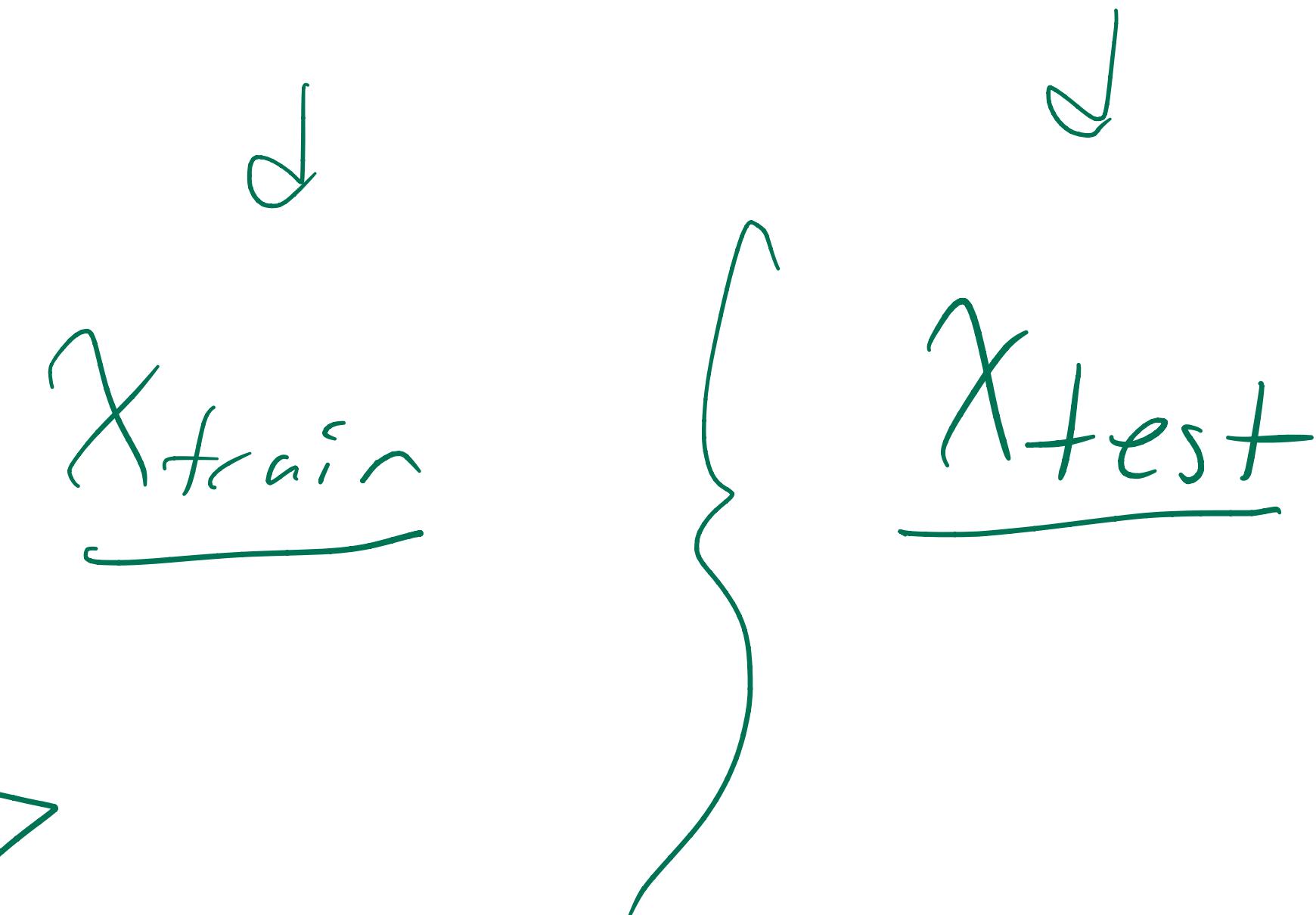
Updating Weights

- if ($f(x_i, w) < 0$ AND $y_i = \text{true}$):
 - $w_0 \leftarrow w_0 + 1$
 - $w_j \leftarrow w_j + \phi(x_i)$
- if ($f(x_i, w) \geq 0$ AND $y_i = \text{false}$):
 - $w_0 \leftarrow w_0 - 1$
 - $w_j \leftarrow w_j - \phi(x_i)$

R

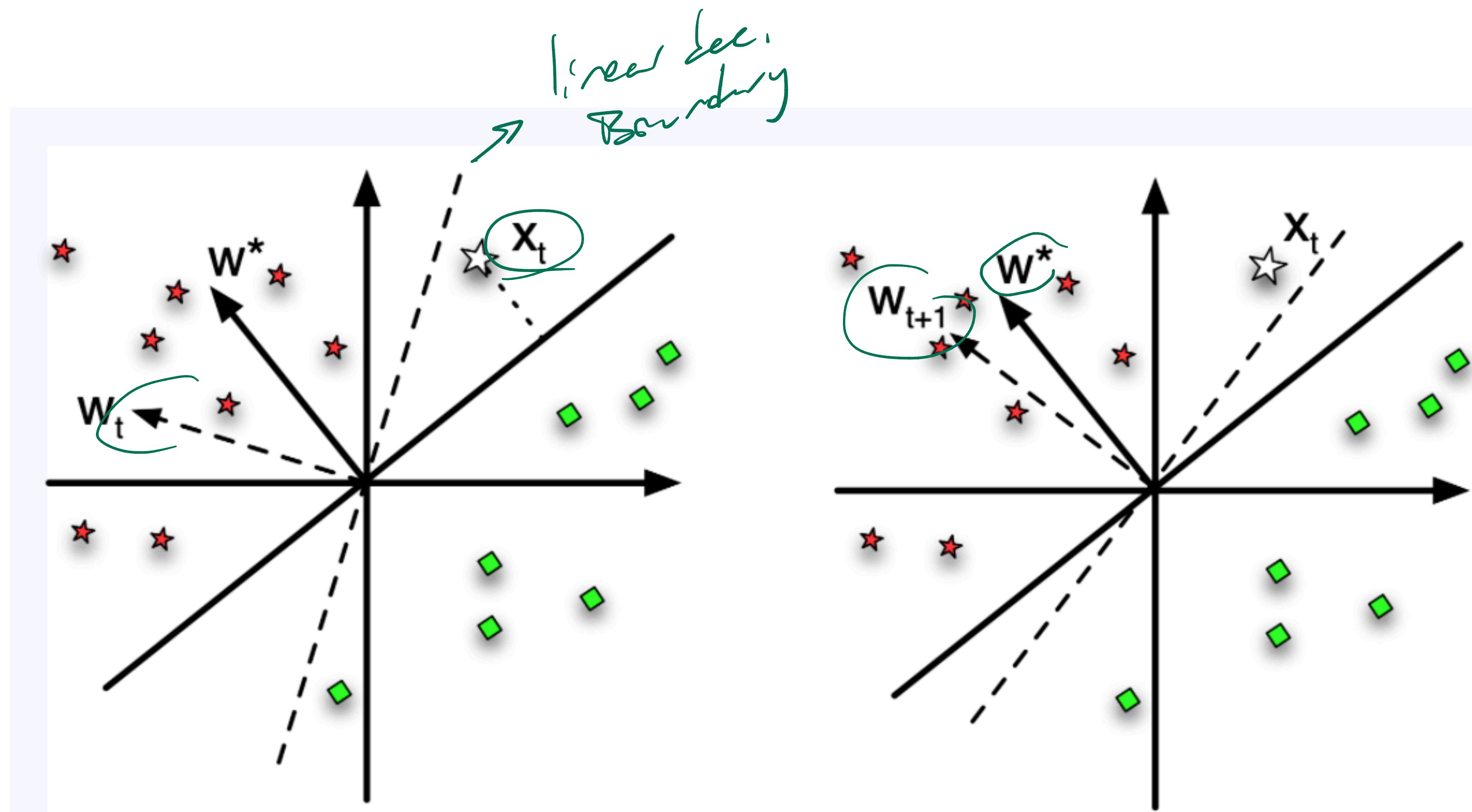
Perceptron Algorithm

1. Initialize weights
2. $\forall (x_i, y_i) \in \mathcal{X}_{train}:$
 - $f(x_i, w) = w_0 + w_1\phi_1(x_1) + \dots + w_n\phi_n(x_n)$
 - If $(f(x_i, w) \geq 0 \text{ AND } y_i = \text{true}) \text{ OR } (f(x_i, w) < 0 \text{ AND } y_i = \text{false}):$
 - Nothing... move to next example
 - Else:
 - if $(f(x_i, w) < 0 \text{ AND } y_i = \text{true}):$
 - $w_0 \leftarrow w_0 + 1$
 - $w_j \leftarrow w_j + \phi(x_i)$
 - if $(f(x_i, w) \geq 0 \text{ AND } y_i = \text{false}):$
 - $w_0 \leftarrow w_0 - 1$
 - $w_j \leftarrow w_j - \phi(x_i)$
3. Repeat until:
 1. No changes
 2. Time limit reached



R

Perceptron

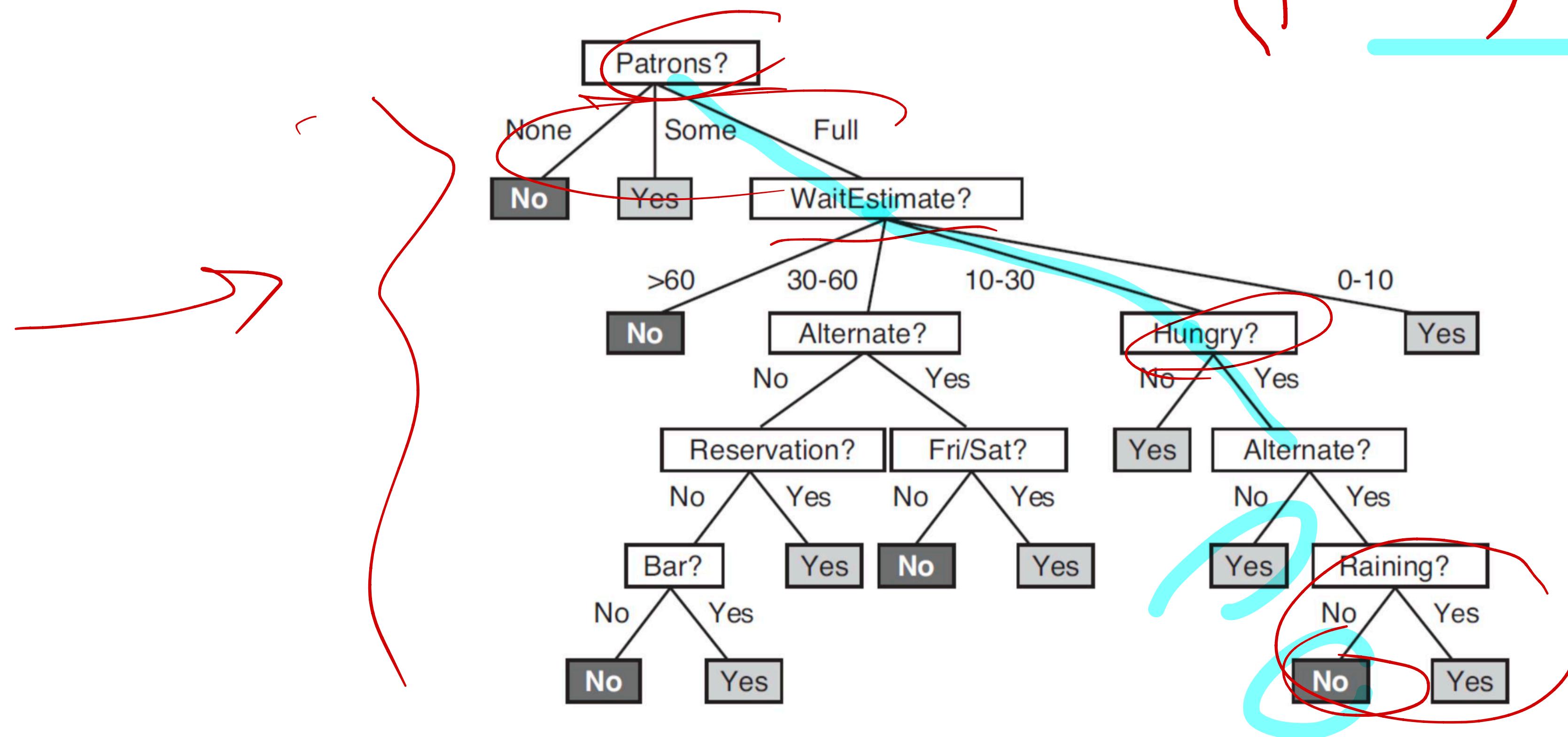


R

Decision Tree

- Decide the outcome based on a set of attributes

(P=Full, WE=20, H=Yes, [-])



R

Decision Tree Algorithm

function LEARN-DECISION-TREE(*examples*, *attributes*, *parent-examples*) **returns** a tree

} if *examples* is empty **then return** PLURALITY-VALUE(*parent-examples*)
else if all *examples* have the same classification **then return** the classification
else if *attributes* is empty **then return** PLURALITY-VALUE(*examples*)
else
 A $\leftarrow \operatorname{argmax}_{a \in \text{attributes}} \text{IMPORTANCE}(a, \text{examples})$
 tree \leftarrow a new decision tree with root test *A*
 for each value *v* of *A* **do**
 exs $\leftarrow \{e : e \in \text{examples} \text{ and } e.A = v\}$
 subtree \leftarrow LEARN-DECISION-TREE(*exs*, *attributes* - *A*, *examples*)
 add a branch to *tree* with label (*A* = *v*) and subtree *subtree*
return *tree*

R

Decision Tree

Algorithm: Importance

- Decide over which attribute to split

- Entropy:

$$H(x) = - \sum_{i=1}^k P(x_i) \log_2 P(x_i)$$

$$B(q) = - (q \log_2(q) + (1-q) \log_2(1-q))$$

$\xrightarrow{\text{To RV with prob } q}$

More entropy \rightarrow Less info

$x \sim \text{fair coin}, x' \sim \text{coin that mostly T}$

$$H(x) = - (.5 \log_2 .5 + .5 \log_2 .5) = 1$$

$$H(x') = - (.99 \log_2 .99 + .01 \log_2 .01) = .066$$

R

Decision Tree

Algorithm: Importance

- Given attribute A with d distinct values: a_1, a_2, \dots, a_d $\rightarrow n^+ \sim \text{Positive labels}$
 $n^- \sim \text{Negative labels}$

- \underline{x}_{train} is divided in d subsets

$$a_k \sim n_k^+ \sim \text{Pos. Labels}$$

$$n_k^- \sim \text{Neg. Labels}$$

$$\text{Remainder}(A) = \sum_{k=1}^d \frac{\frac{n_k^+ + n_k^-}{n^+ + n^-} B - \frac{n_k^+}{n_k^+ + n_k^-}}{\text{total}}$$

R

Decision Tree

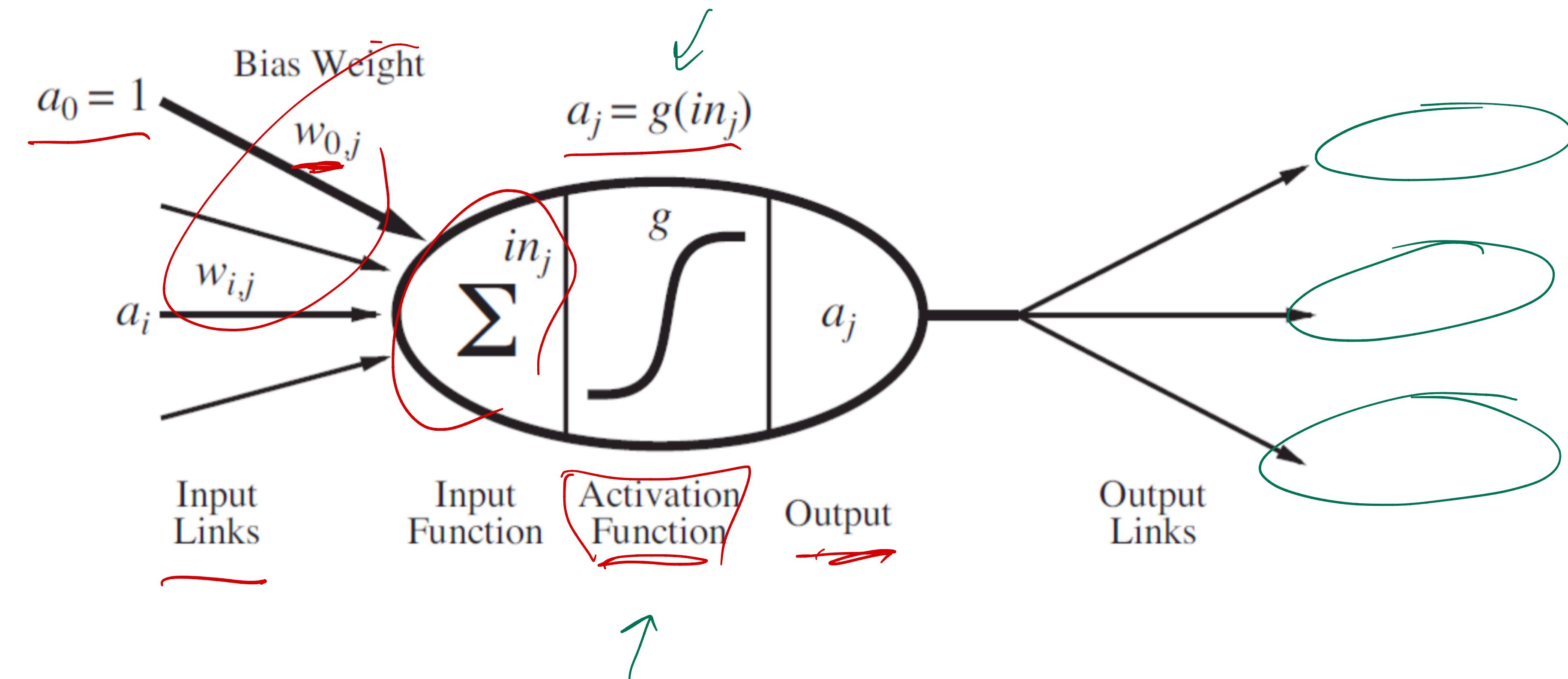
Algorithm: Importance

- Information gain: the expected reduction in entropy

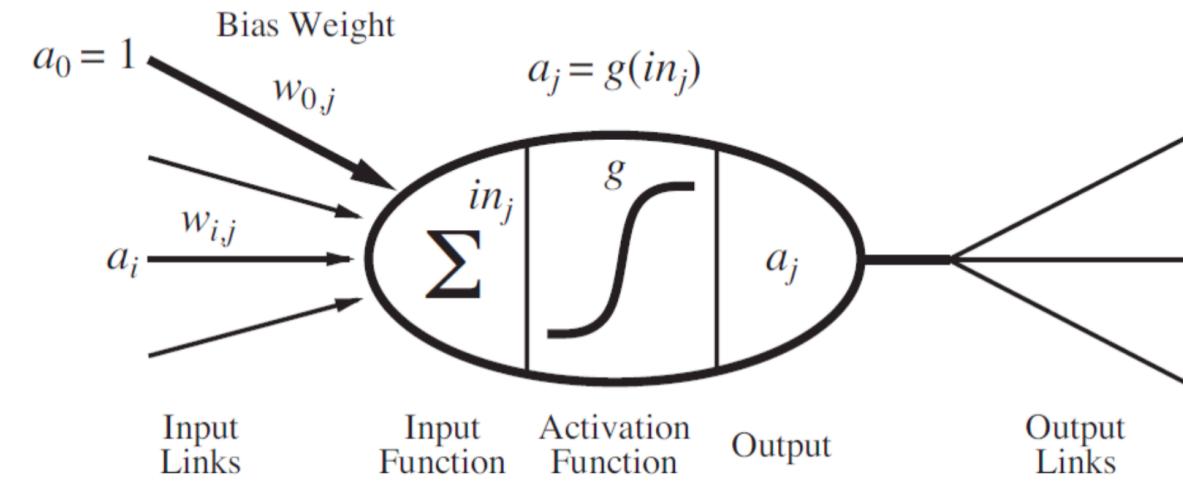
$$\text{Importance} = \underline{\text{Gain}}(A) = B\left(\frac{n^+}{n^+ + n^-}\right) - \underline{\text{Remainder}}(A)$$

R

Neuron Mathematical Model

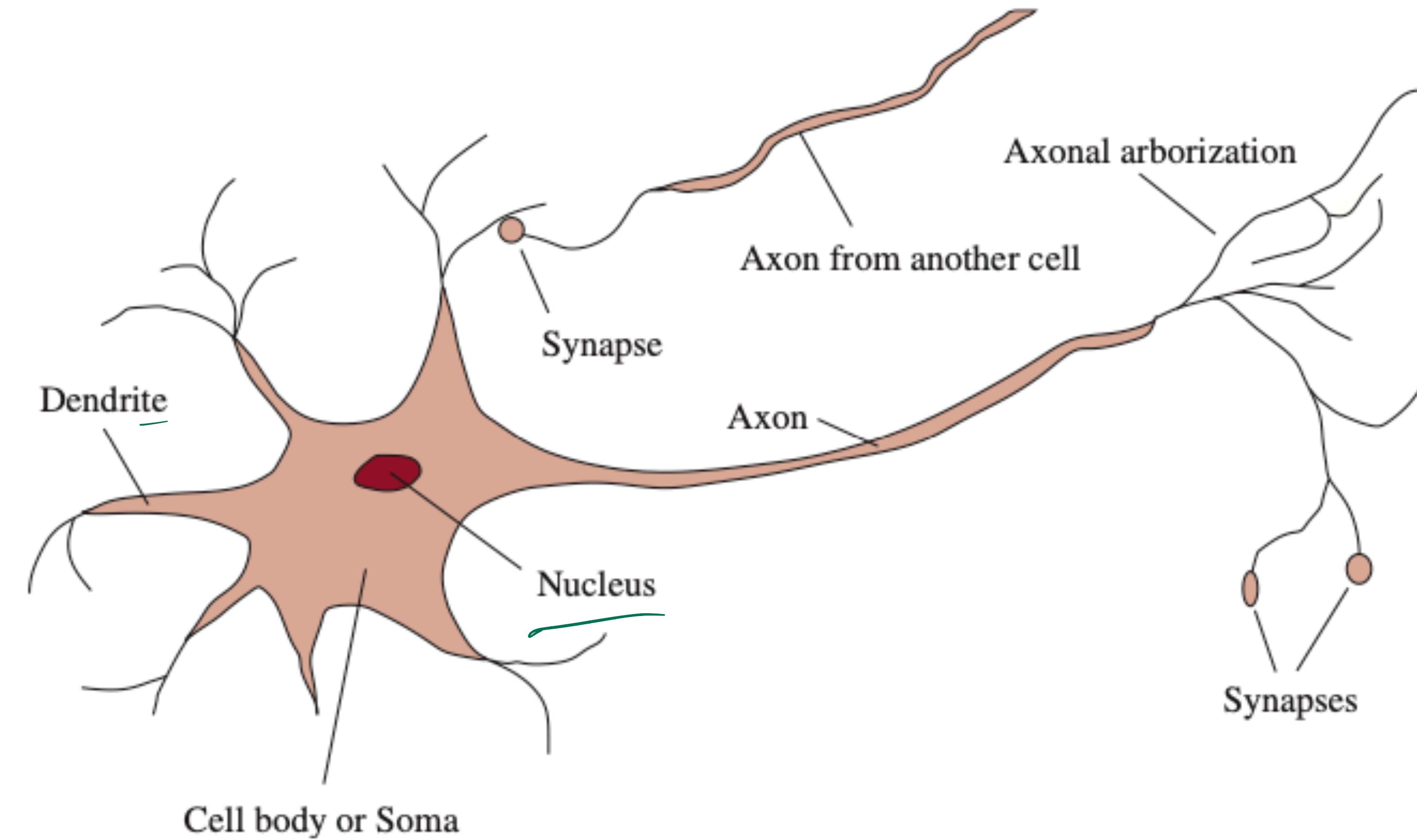


R



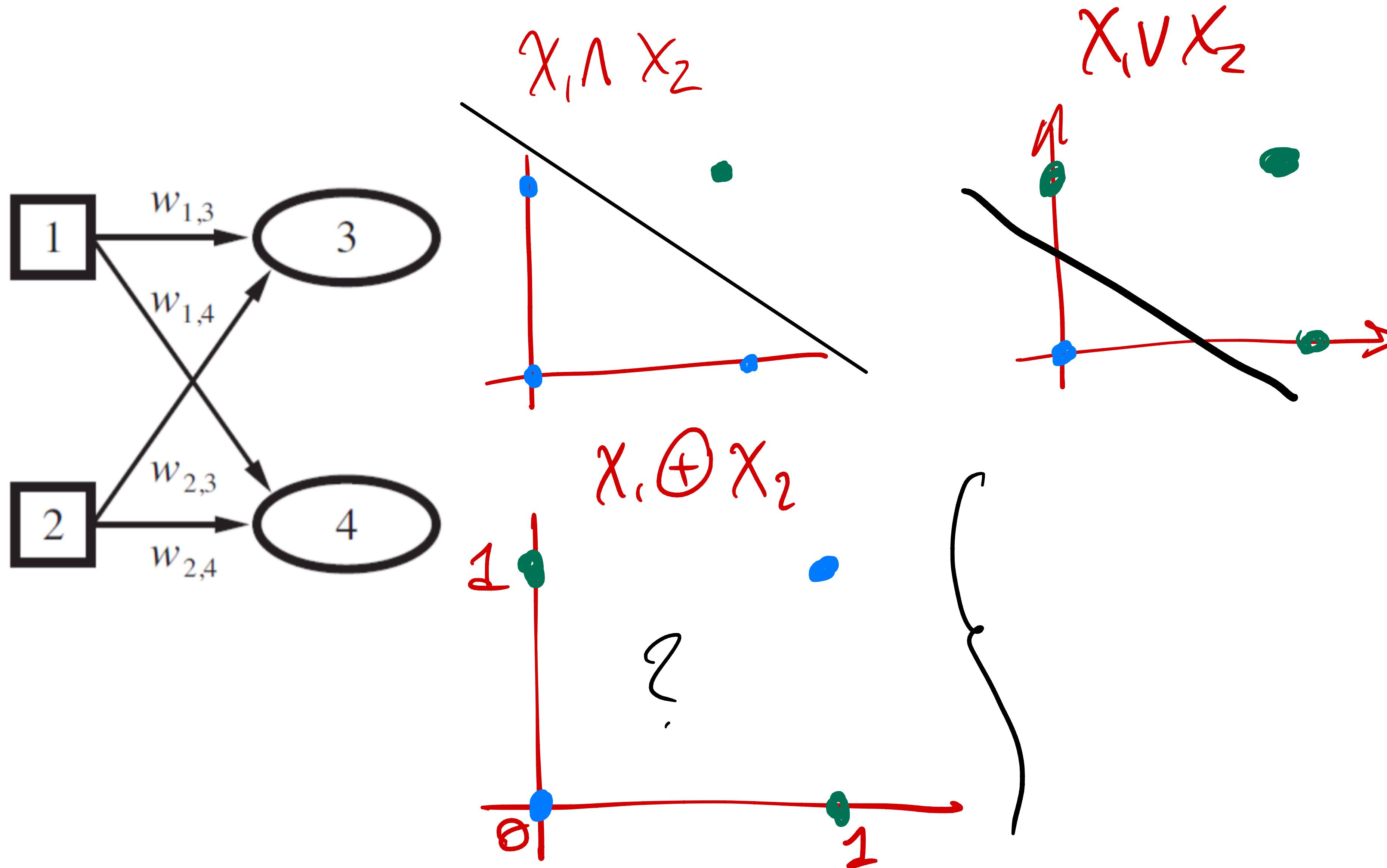
Neuron

Real Neuron



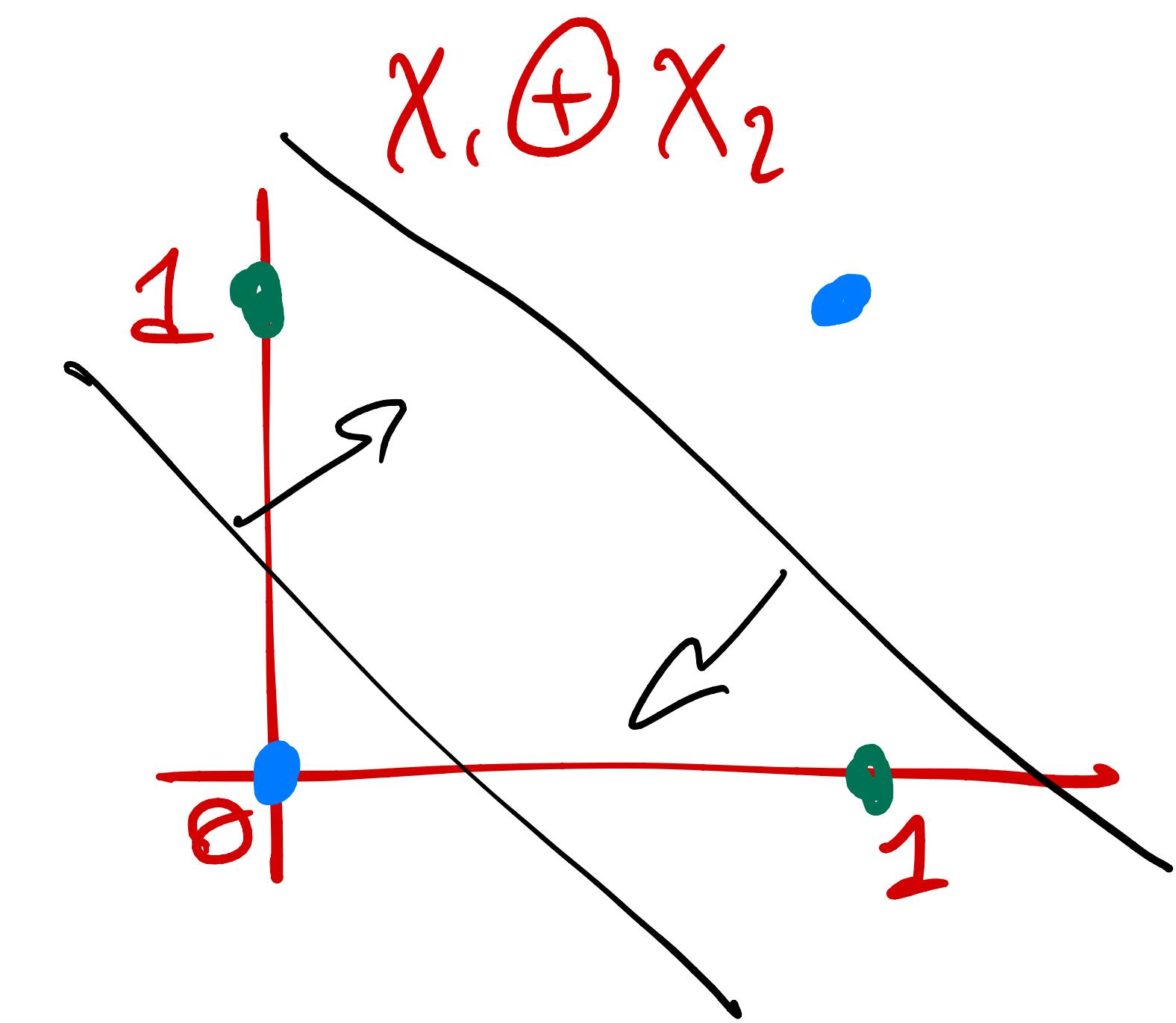
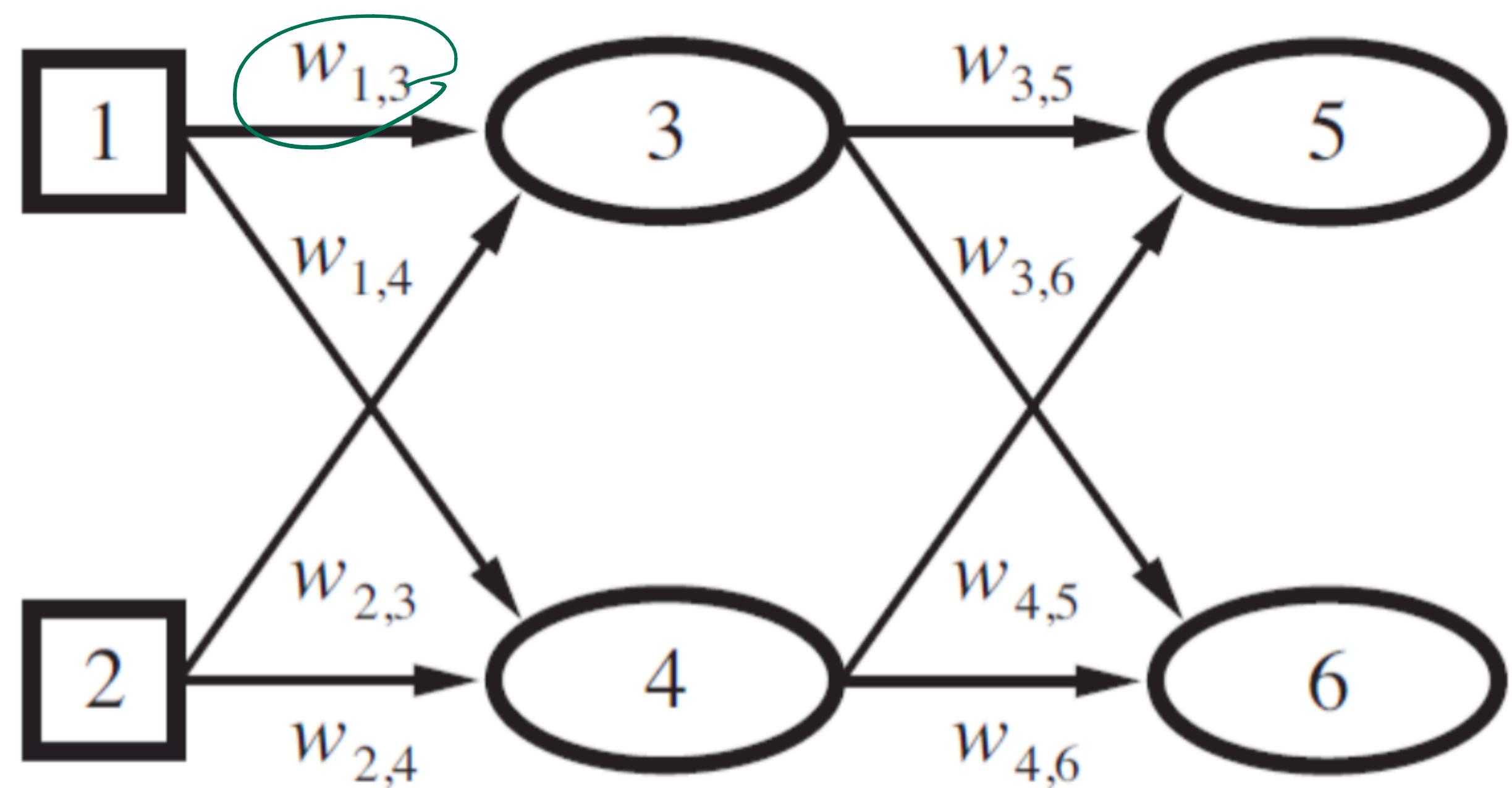
R

Single Layer Feed-Forward NN



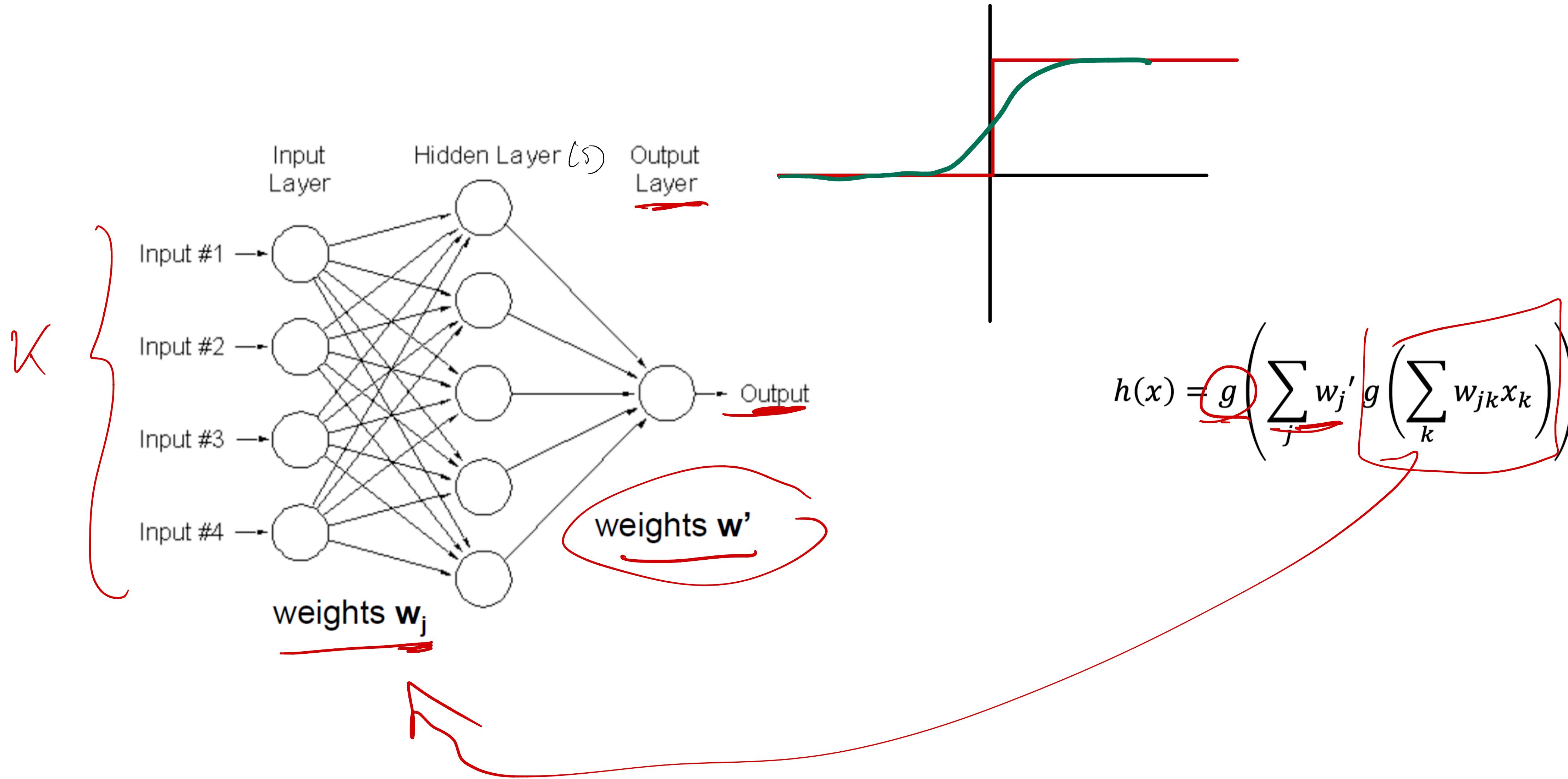
R

Multi layer



R

General Structure

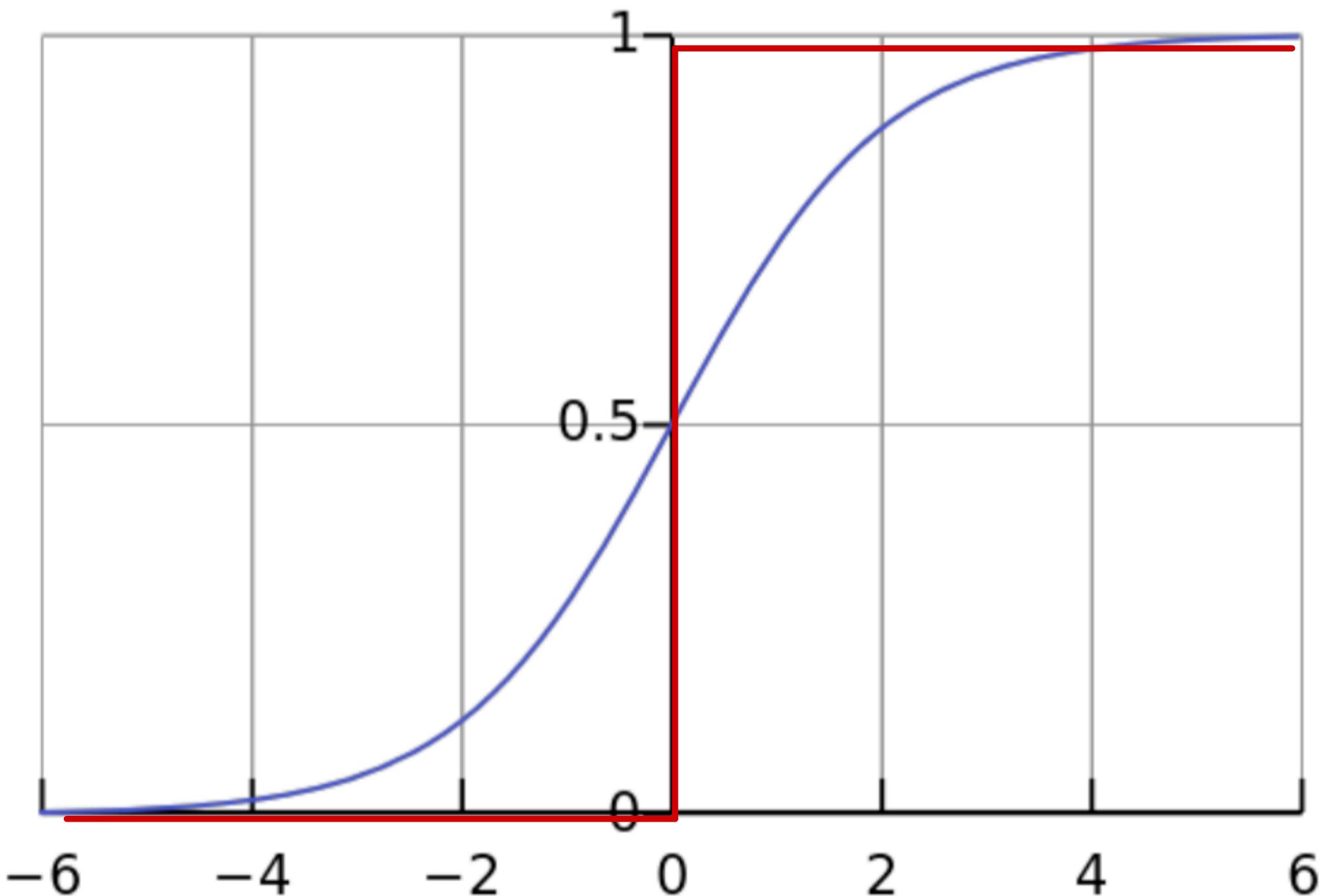


R

Activation function

- Sigmoid:

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

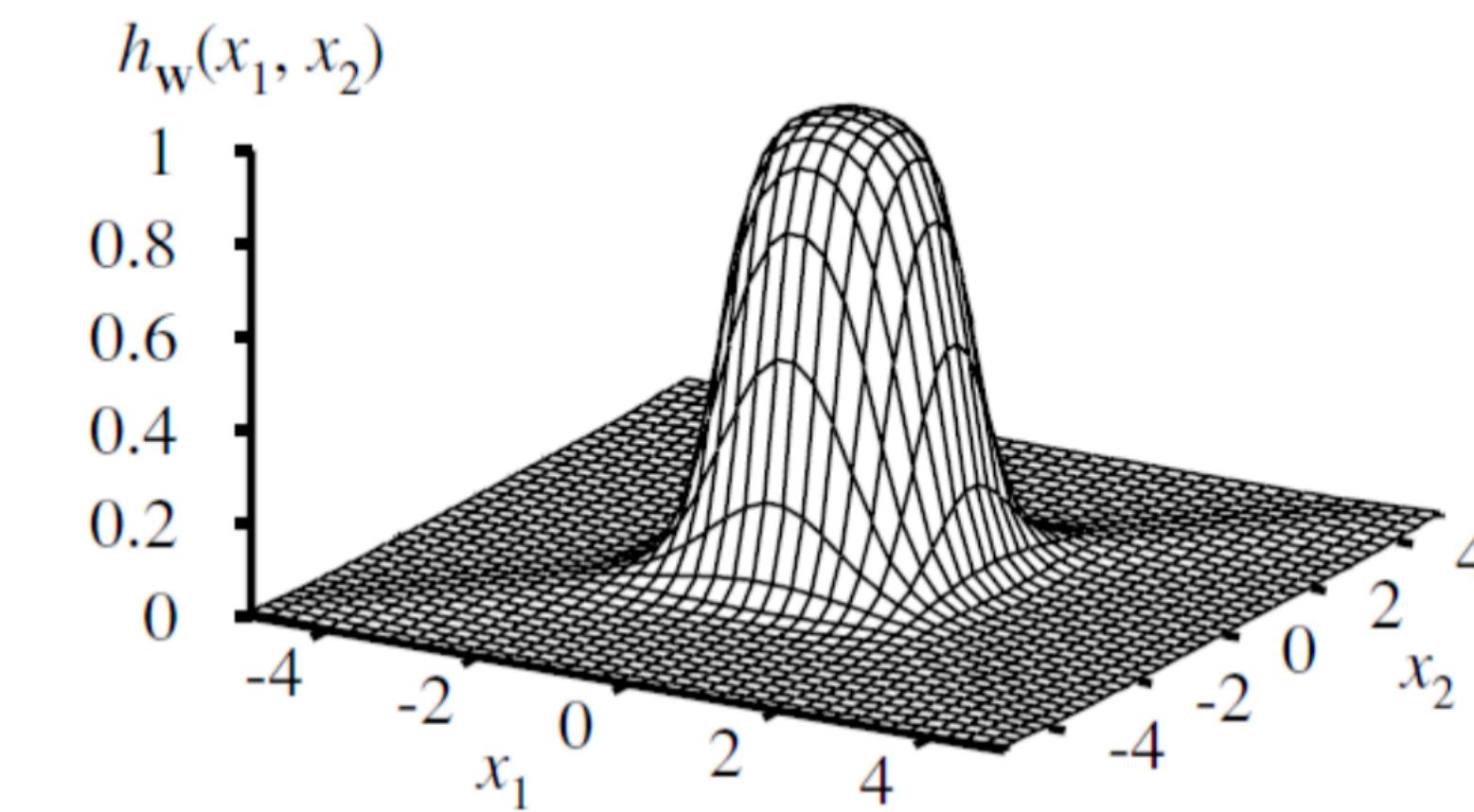
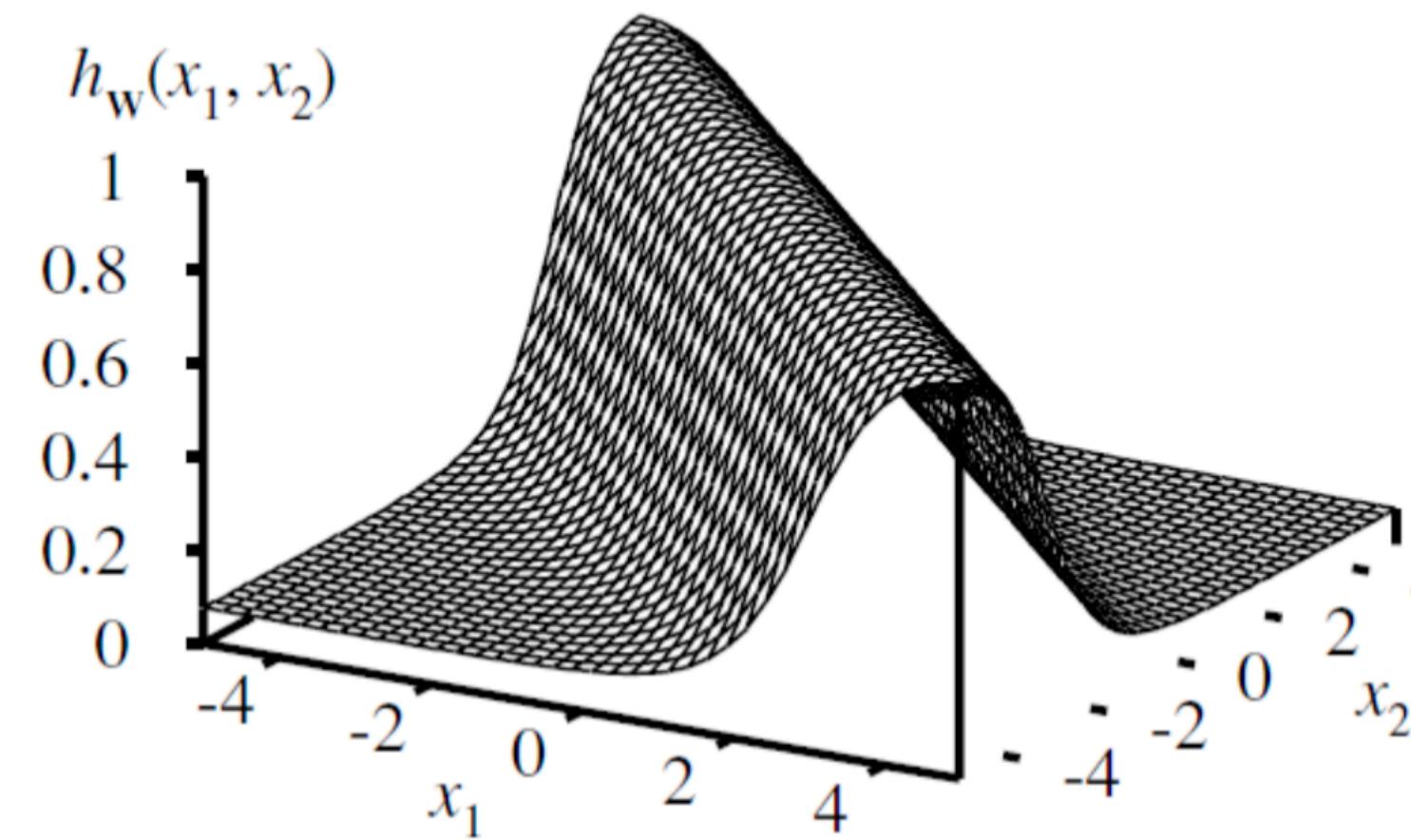


R

Activation function

- Sigmoid:

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



R

Differentiable Perceptron

- Define the error

$$E(w) = -\frac{1}{2} \sum_j \underline{(y_j - h_w(x_j))^2}, \quad h_w(x_j) = g(w^T x_j)$$

- Update via gradient descent

$$w \leftarrow w - \alpha \frac{\partial E}{\partial w}$$

$$w \leftarrow w - \alpha [y_i - g(w^T x_j)] [g(w^T x_j)(1 - g(w^T x_j)) x_j]$$

$$w \leftarrow w + \alpha (y_i - h_w(x_j)) x_j$$

R

Artificial Neural Networks

Comparing

PROS

- Easy to implement
(TensorFlow)
- Approximate Functions
- More layers → Better Results

CONS

- Depend on X
 - Stuck in local optima
 - Difficult to analyze
-

R

Deep Neural Networks

- end to end training
- Experience needed for Net Structure
- Requires
 - ↑ Data
 - ↑ Computing

