

Einführung in Artificial Intelligence SS 2024, 4.0 VU, 192.027

Exercise Sheet 3 – Learning from Examples and Neural Networks

For the discussion part of this exercise, mark and upload your solved exercises in **TUWEL** until Wednesday, June 5, 23:55 CEST. The registration for a solution discussion ends on Friday, June 7, 23:55 CEST. Be sure that you tick only those exercises that you can solve and explain!

In the discussion, students will be asked questions about their solutions of examples they checked. The discussion will be evaluated with 0–25 points, which are weighted with the fraction of checked examples and rounded to the next integer. There is *no minimum number of points* needed for a positive grade (i.e., you do not need to participate for a positive grade, but you can get at most $\approx 80\%$ without doing exercises).

Note, however, that *your registration is binding*. Thus, *if* you register for a solution discussion, then it is *mandatory* to show up. No-show after a registration without plausible excuse leads to a penalty. Such students will not have the possibility to participate in another exam once they have gotten a certificate. If you registered but cannot show up due to unpredictable and unavoidable obstacles, either deregister or send us a confirmation (doctor's note, etc) and you will be excused. Please ask questions in the **TUWEL** forum or visit our tutors during the tutor hours (see **TUWEL**).

Exercise 3.1: A close friend of you is experiencing some trouble in deciding which Lego set to buy next. You have recently learned about the wonders of artificial intelligence. In particular, how decision trees can be used to make predictions based on examples and you decide to construct a model that decides whether your friend will like a newly released Lego set or not. After a quick research about the relevant data for such a decision, you settle for the attributes P (the number of pieces in the set) with the domain $V(P) = \{<500, 500-1000, >1000\}$, R (whether the release year is after 2010) with the domain $V(R) = \{\top, \perp\}$, S (whether the set contains a lot of stickers) with the domain $V(S) = \{\top, \perp\}$ and T (theme the set belongs to) with the domain $V(T) = \{\text{star wars, superheros, technic, architecture, city, creator, ninjago}\}$. Content with this rather simplistic model, your friend provides you with the data for some sets that they have already build:

Sample	P	R	S	T	Liked?
1	500–1000	\top	\perp	star wars	T
2	<500	\perp	\top	ninjago	F
3	>1000	\perp	\top	technic	T
4	>1000	\top	\perp	architecture	T
5	<500	\perp	\perp	creator	T
6	500–1000	\top	\top	superheros	F
7	500–1000	\perp	\top	city	F
8	500–1000	\top	\perp	technic	T
9	>1000	\perp	\top	star wars	T

Use the gathered data to construct a decision tree capable of predicting whether your friend will like a set, given the values of the attributes they chose. In each step of the construction, choose the attribute that maximizes the information gain, as shown in the lecture.

Exercise 3.2: Construct another decision tree for Exercise 3.1, this time choosing the attribute that maximizes the *relative information gain*, i.e., the ratio between the gain of the attribute and its own intrinsic information.

$$\text{GainR}(A) = \frac{\text{Gain}(A)}{H(A)}$$

and

$$H(A) := \sum_{a \in V(A)} \frac{|E_a|}{S} \log_2 \frac{S}{|E_a|},$$

where $V(A)$ denotes the domain of the attribute A and $|E_a|$ is the number of all samples which have the value a of the attribute A . Furthermore, we define $S := \sum_{a \in V(A)} |E_a|$ to be the size of the set of all samples.

Exercise 3.3: In this exercise we explore some of the problems with decision trees and the ways of dealing with them:

- (a) Compare the results obtained in Exercise 3.1 and Exercise 3.2. Use this example to discuss the advantages of using the relative information gain rather than the regular information gain rule. Be sure to explain *why* using the relative gain leads to better results!
- (b) Discuss the design choices that were made in Exercise 3.1.

Are decision trees a suitable model for predicting which Lego sets someone will like *in practice*?

Is the choice of attributes and possible values sensible? If so, argue why, if not, provide some alternative options.

What can you say about the practical accuracy of the generated tree(s)? What could have been done to make the tree(s) more accurate (match your friend's situation better)?

- (c) Suppose that you are collecting data for a decision tree. The first two examples you collect have the exact same value for each of the attributes you picked, but their classification is different. You stop collecting further data and ponder the implications of this situation. Answer the following questions briefly:
 - What could be the cause of such a situation?
 - What would the decision tree learning algorithm discussed in the lecture do in this case?
 - What could you do to avoid the problem?

Exercise 3.4: Suppose that an attribute splits the set of examples E into k subsets E_1, \dots, E_k where each subset E_i has p_i positive and n_i negative examples. Show that the attribute has zero information gain if the ratio $\frac{p_i}{p_i+n_i}$ is equal for all $i \in \{1, \dots, k\}$. Provide arguments for all properties that you use.

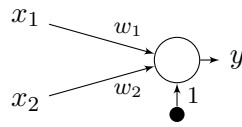
Exercise 3.5: Construct a model of a neural network with three binary input signals I_1, I_2, C_{in} and two output signals C_{out} (carry), S (sum) which represents an implementation of a *full adder*. The signals should behave according to the entries of the following truth table.

I_1	I_2	C_{in}	C_{out}	S
0	0	0	0	0
0	0	1	0	1
0	1	0	0	1
0	1	1	1	0
1	0	0	0	1
1	0	1	1	0
1	1	0	1	0
1	1	1	1	1

For each neuron, use the following activation function:

$$g(x) = \begin{cases} 1 & \text{if } x \geq 0.5, \\ 0 & \text{otherwise.} \end{cases}$$

Exercise 3.6: Consider a single layer perceptron with two input neurons and one output neuron of the following form:



Train the perceptron using the *Perceptron learning rule* and the identity activation function $g(x) = x$ on the following training data: $f(1, 2) = 2$, $f(2, 1) = 7$, $f(2, -1) = 4$ and $f(8, 2) = 7$. The weights are initialized to 0 and the learning rate α is 1. The bias weight will not be learned and stays 1.

Express the function $f(x_1, x_2)$ that your network has learned after the training is complete algebraically. Does this function produce the expected outputs for all four training inputs?

If the neuron has not learned the function correctly, provide alternative weights so that the neuron produces the expected outputs for all examples, or prove that this is impossible.

Exercise 3.7: Show formally that a single-layer perceptron with the step function as activation function ($g(x) = 1$ for $x \geq t$ for some threshold t and $g(x) = 0$ otherwise) cannot express the logical equivalence operator (\equiv).

Exercise 3.8:

1. Explain from a *gradient descent* viewpoint that in the multi-variable linear regression model, L_1 regularization prefers to set many weights to 0 which results in sparse models, comparing with L_2 regularization.

Recall that the L_q regularization adds to the (empirical) loss function a regularization term $\lambda \cdot \text{Complexity}(h_{\mathbf{w}})$, where λ is a number and $\text{Complexity}(h_{\mathbf{w}}) = \sum_i |w_i|^q$, with $q \in \{1, 2\}$.

Hint: observe the behavior of function *signum*, the derivative of absolute value function.

2. Show that the sigmoid function σ is a special softmax function with $d = 2$.

Recall the two functions as follows.

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

$$\text{softmax}(\mathbf{v})_k = \frac{e^{v_k}}{\sum_{k'}^d e^{v_{k'}}}, \mathbf{v} = v_1 \dots v_d.$$

3.1

Attributes A:

$$P \in V(P) = \{ < 500, 500-1000, > 1000 \}$$

$$R \in V(R) = \{ T, \perp \}$$

$$S \in V(S) = \{ T, \perp \}$$

$$T \in V(T) = \{ \text{star wars, superheroes, technic, architecture, city, creator, ninjago} \}$$

↳ choose Attribute A with maximal $\text{Gain}(A)$

$$\text{Gain}(A) = B\left(\frac{p}{p+n}\right) - \text{Rem}(A)$$

$$\text{Rem}(A) = \sum_{k=1}^d \frac{p_k + n_k}{p + n} B\left(\frac{p_k}{p_k + n_k}\right)$$

$\text{Rem}(A)$... Remainder

d ... number of distinct values in the domain/set $V(A)$

E_k ... Attribute A with d distinct values divides into subsets E_1, E_2, \dots, E_d

p ... number of total positive examples

n ... number of total negative examples

p_k ... number of positive examples for subset E_k

n_k ... number of negative examples for subset E_k

$$B(q) = - \left[q \log_2(q) + (1-q) \log_2(1-q) \right]$$

$B(q)$... Entropy of a boolean random variable that is true with probability q .

① Generic evaluations

$$p=6 \quad n=3 \quad \text{samples}=9$$

$$B\left(\frac{p}{p+n}\right) = B\left(\frac{6}{9}\right) = - \left[\frac{1}{3} \log_2\left(\frac{1}{3}\right) + \left(1 - \frac{1}{3}\right) \log_2\left(1 - \frac{1}{3}\right) \right]$$

$$B\left(\frac{1}{3}\right) = \frac{\log_2\left(\frac{27}{4}\right)}{\log_2(8)} \approx 0,918$$

② choose first attribute:

Attribute P:

$$E_{<500} = \left\{ \begin{array}{l} p_1 = 1 \\ n_1 = 7 \\ |E_{<500}| = 2 \end{array} \right\} \quad E_{500-1000} = \left\{ \begin{array}{l} p_2 = 2 \\ n_2 = 2 \\ |E_{500-1000}| = 4 \end{array} \right\} \quad E_{>1000} = \left\{ \begin{array}{l} p_3 = 3 \\ n_3 = 0 \\ |E_{>1000}| = 3 \end{array} \right\}$$

$$Rem(P) = \underbrace{\frac{2}{9} B\left(\frac{1}{2}\right)}_{E_{<500}} + \underbrace{\frac{4}{9} B\left(\frac{2}{4}\right)}_{E_{500-1000}} + \underbrace{\frac{3}{9} B\left(\frac{3}{3}\right)}_{E_{>1000}}$$

$$Rem(P) = \frac{2}{9} + \frac{4}{9} + 0 = \frac{2}{3}$$

$$Gain(P) = B\left(\frac{1}{3}\right) - Rem(P) = 0,918 - \frac{2}{3} = 0,251$$

Attribute R:

$$E_T = \left\{ \begin{array}{l} p=3 \\ n=1 \\ |E_T|=4 \end{array} \right\} \quad E_L = \left\{ \begin{array}{l} p=3 \\ n=2 \\ |E_L|=5 \end{array} \right\}$$

$$Rem(R) = \frac{4}{9} B\left(\frac{3}{4}\right) + \frac{5}{9} B\left(\frac{3}{5}\right) = 0,361 + 0,539 = 0,9$$

$$Gain(R) = 0,918 - 0,9 = 0,018$$

Attribute S:

$$E_T = \left\{ \begin{array}{l} p=2 \\ n=3 \\ |E_T|=5 \end{array} \right\} \quad E_L = \left\{ \begin{array}{l} p=4 \\ n=0 \\ |E_L|=4 \end{array} \right\}$$

$$Rem(S) = \frac{5}{9} B\left(\frac{2}{5}\right) + \frac{4}{9} B\left(\frac{4}{4}\right) = 0,539$$

$$Gain(S) = 0,918 - 0,539 = 0,379$$

Attribute T:

$$E_{star\ wars} = \left\{ \begin{array}{l} p=2 \\ n=0 \\ |E_1|=2 \end{array} \right\} \quad E_{ninjaigo} = \left\{ \begin{array}{l} p=0 \\ n=1 \\ |E_2|=1 \end{array} \right\} \quad E_{technik} = \left\{ \begin{array}{l} p=2 \\ n=0 \\ |E_3|=2 \end{array} \right\}$$

$$E_{architecture} = \left\{ \begin{array}{l} p=1 \\ n=0 \\ |E_4|=1 \end{array} \right\} \quad E_{creator} = \left\{ \begin{array}{l} p=1 \\ n=0 \\ |E_5|=1 \end{array} \right\} \quad E_{superheros} = \left\{ \begin{array}{l} p=0 \\ n=1 \\ |E_6|=1 \end{array} \right\}$$

$$E_{city} = \left\{ \begin{array}{l} p=0 \\ n=1 \\ |E_7|=1 \end{array} \right\}$$

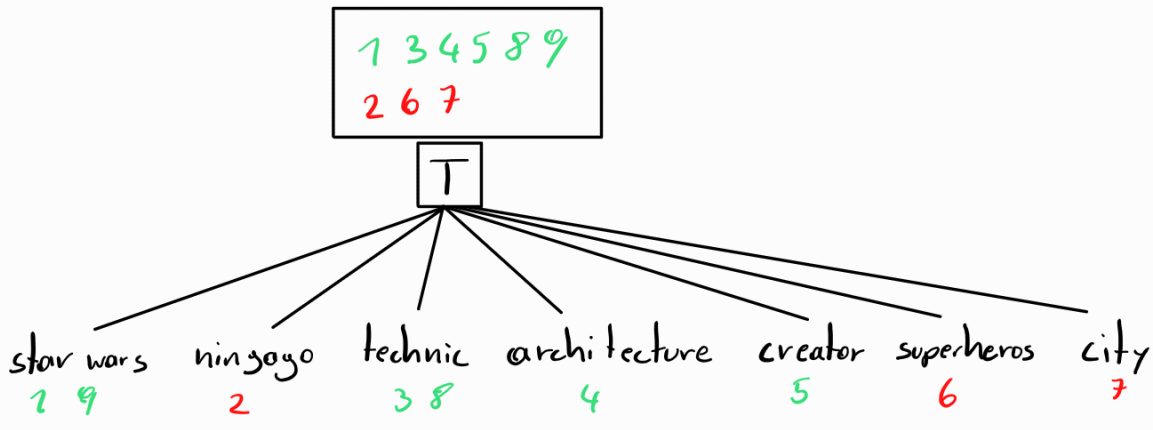
$$Rem(T) = \underbrace{\frac{2}{9} B(\frac{2}{2})}_0 + \underbrace{\frac{1}{9} B(\frac{1}{1})}_0 + \underbrace{\frac{2}{9} B(\frac{2}{2})}_0 + \underbrace{\frac{1}{9} B(\frac{1}{1})}_0 + \underbrace{\frac{1}{9} B(\frac{1}{1})}_0 + \underbrace{\frac{1}{9} B(\frac{0}{1})}_0 + \underbrace{\frac{1}{9} B(\frac{0}{1})}_0$$

$$Rem(T) = 0$$

$$Gain(T) = 0,978 - 0 = 0,978$$

$$Rem(T) > Rem(S) > Rem(P) > Rem(R)$$

First Attribute = T



3.2

$$\text{Gain}_R(A) = \frac{\text{Gain}(A)}{H(A)}$$

$$H(A) = \sum_{\alpha \in V(A)} \frac{|E_\alpha|}{S} \log_2 \frac{S}{|E_\alpha|}$$

$$S = \sum_{\alpha \in V(\alpha)} |E_\alpha| = p + k = \text{number of all samples}$$

① Choose first Attribute

Attribute P:

$$\text{Gain}(P) = 0,257 \quad S = 9$$

$$H(P) = \frac{2}{9} \log_2 \frac{9}{2} + \frac{4}{9} \log_2 \frac{9}{4} + \frac{3}{9} \log_2 \frac{9}{3}$$

$$H(P) = 0,482 + 0,52 + 0,528 = 1,53$$

$$\text{Gain}_R(P) = \frac{0,084}{1,53} = 0,164$$

Attribute R:

$$\text{Gain}(R) = 0,018 \quad S = 9$$

$$H(R) = \frac{4}{9} \log_2 \frac{9}{4} + \frac{5}{9} \log_2 \frac{9}{5}$$

$$H(R) = 0,52 + 0,471 = 0,991$$

$$\text{Gain}_R(H) = \frac{0,018}{0,991} = 0,018$$

Attribute S:

$$\text{Gain}(S) = 0,379 \quad S = 9$$

$$H(S) = \frac{5}{9} \log_2 \frac{9}{5} + \frac{4}{9} \log_2 \frac{9}{4} = 0,991$$

$$\text{Gain}_R(S) = \frac{0,379}{0,991} = 0,382$$

Attribute T:

$$Gain(T) = 0,978$$

$$H(T) = 2 \cdot \frac{2}{9} \log_2 \frac{9}{2} + 5 \cdot \frac{1}{9} \log_2 \frac{9}{1}$$

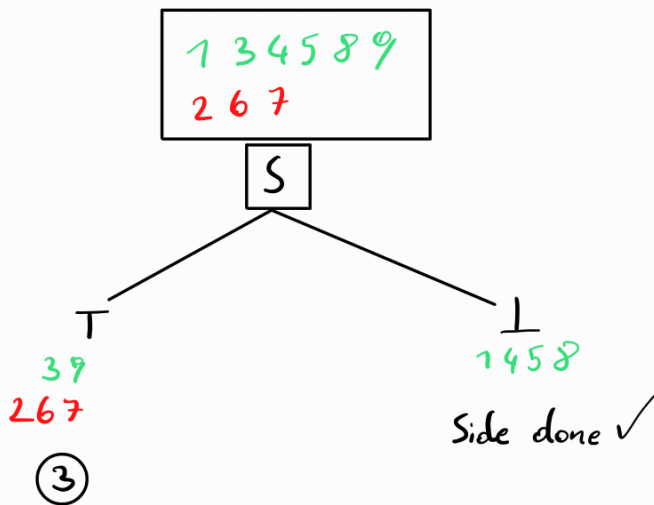
$$H(T) = 2 \cdot \frac{\log_2(\frac{81}{4})}{\log_2(572)} + 5 \cdot \frac{\log_2(9)}{\log_2(572)}$$

$$H(T) = 0,964 + 1,761 = 2,725$$

$$GainR(T) = \frac{0,978}{2,725} = 0,337$$

First Attribute = S

② Construct partial tree



③ Best attribute for branch S=T

Consider only samples 2, 3, 6, 7, 9

total samples = 5 p = 2 n = 3

$$B\left(\frac{p}{p+n}\right) = B\left(\frac{2}{5}\right) = - \left[\frac{2}{5} \log_2\left(\frac{2}{5}\right) + \left(1 - \frac{2}{5}\right) \log_2\left(1 - \frac{2}{5}\right) \right]$$

$$B\left(\frac{1}{2}\right) = 0,977$$

Attribute P:

$$E_{>1000} = \left\{ \begin{array}{l} p=2 \\ n=0 \\ |E|=2 \end{array} \right\} \quad E_{500-1000} = \left\{ \begin{array}{l} p=0 \\ n=2 \\ |E|=2 \end{array} \right\} \quad E_{<500} = \left\{ \begin{array}{l} p=0 \\ n=1 \\ |E|=1 \end{array} \right\}$$

$$Rem(P) = \frac{2}{4} \underbrace{B\left(\frac{2}{2}\right)}_0 + \frac{2}{4} \underbrace{B\left(\frac{0}{2}\right)}_0 + \frac{1}{4} B\left(\frac{0}{1}\right) = 0$$

$$Gain(P) = 0,977 - 0 = 0,977$$

$$H(P) = 2 \cdot \frac{2}{5} \log_2\left(\frac{5}{2}\right) + \frac{1}{5} \log_2\left(\frac{5}{1}\right) = 1,522$$

$$GainR(P) = \frac{Gain(H)}{H(H)} = \frac{0,977}{1,522} = 0,634$$

Attribute R:

$$E_T = \left\{ \begin{array}{l} p=0 \\ n=1 \\ |E_T|=1 \end{array} \right\} \quad E_L = \left\{ \begin{array}{l} p=2 \\ n=2 \\ |E_L|=4 \end{array} \right\} \quad S=5$$

$$Rem(R) = \underbrace{\frac{1}{5} B\left(\frac{0}{1}\right)}_0 + \underbrace{\frac{4}{5} B\left(\frac{2}{4}\right)}_{\frac{4}{5}} = \frac{4}{5}$$

$$Gain(R) = 0,9777 - \frac{4}{5} = 0,7777$$

$$H(R) = \frac{1}{5} \log_2\left(\frac{5}{1}\right) + \frac{4}{5} \log_2\left(\frac{5}{4}\right) = 0,7222$$

$$GainR(R) = \frac{0,7777}{0,7222} = 0,237$$

Attribute T:

$$E_{technic} = \left\{ \begin{array}{l} p=1 \\ n=0 \\ |E|=1 \end{array} \right\} \quad E_{superheros} = \left\{ \begin{array}{l} p=0 \\ n=1 \\ |E|=1 \end{array} \right\} \quad E_{city} = \left\{ \begin{array}{l} p=0 \\ n=1 \\ |E|=1 \end{array} \right\} \quad E_{star\ vars} = \left\{ \begin{array}{l} p=1 \\ n=0 \\ |E|=1 \end{array} \right\} \quad E_{ninjas} = \left\{ \begin{array}{l} p=0 \\ n=1 \\ |E|=1 \end{array} \right\}$$

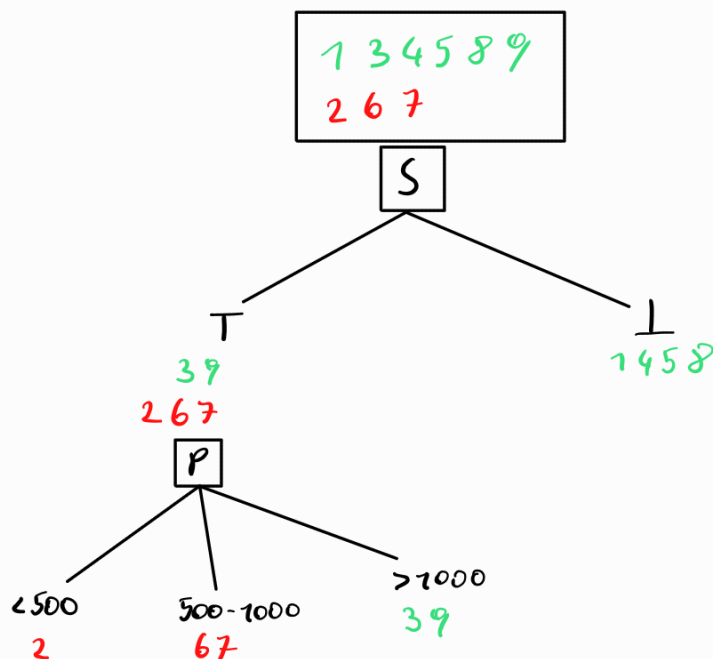
$$Rem(T) = \frac{1}{5} B\left(\frac{1}{1}\right) + \frac{1}{5} B\left(\frac{0}{1}\right) + \frac{1}{5} B\left(\frac{0}{1}\right) + \frac{1}{5} B\left(\frac{1}{1}\right) + \frac{1}{5} B\left(\frac{0}{1}\right) = 0$$

$$Gain(T) = 0,9777 - 0 = 0,9777$$

$$H(T) = 5 \cdot \left[\frac{1}{5} \cdot \log_2\left(\frac{5}{1}\right) \right] = 2,322$$

$$Gain(T) = \frac{1}{2}$$

Selected Attribute = P



3.3

(a)

By splitting the samples based on the regular information gain, the algorithm will choose the most relevant attribute, in such a fashion that returns as many subsets with (ideally) all positive/negative classifications as possible. This usually means that an attribute with a lot of distinct values is chosen, resulting in a higher branching factor (a wider tree), as seen in 3.7

In order to avoid this, we can use the relative information gain, because by dividing $\frac{\text{gain}(A)}{H(A)}$ the gain by its intrinsic information we take the number of values into account.

(b)

• Are decision trees a suitable model for predicting which Lego sets someone will like in practice?

Yes, if the attributes are chosen appropriately

• Is the choice of attributes and possible values sensible? If so, argue why, if not, provide some alternate options

The attribute S is loosely defined and should be expressed as a numerical value, in order to set thresholds.

• What can you say about the practical accuracy of the generated trees. What could have been done to make the trees more accurate (match your friend's situation better)?

In order to increase the accuracy of the tree, we need to increase the sample size. Because $n=9$ in 3.7, the tree might not be very accurate in practical applications

(c)


Classification: Outcome i.e. Boolean classification: True/False

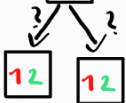
• What could be the cause of such a situation?

A missing additional attribute that would differentiate the two examples.

• What would the decision tree learning algorithm discussed in the lecture do in this case?

The algorithm is unlikely to terminate, because there is no way to differentiate the two samples.

Node  could not be split based on any existing attributes.



• What could you do to avoid the problem?

Add an additional attribute to differentiate the two

3.4

$$Gain(A) = B\left(\frac{P}{P+n}\right) - Rem(A)$$

$$Rem(A) = \sum_{k=1}^d \frac{P_k + n_k}{P + n} B\left(\frac{P_k}{P_k + n_k}\right)$$

$$Gain(A) = B\left(\frac{P}{P+n}\right) - \sum_{k=1}^d \frac{P_k + n_k}{P + n} B\left(\frac{P_k}{P_k + n_k}\right)$$

$$Gain(A) = B\left(\frac{P}{P+n}\right) - B\left(\frac{P_k}{P_k + n_k}\right) \sum_{k=1}^d \frac{P_k + n_k}{P + n}$$

$$Gain(A) = B\left(\frac{P}{P+n}\right) - B\left(\frac{P_k}{P_k + n_k}\right)$$

$$Gain(A) = 0$$

Because $\frac{P_k}{P_k + n_k}$ is the same for every k of the sum we can treat $B\left(\frac{P_k}{P_k + n_k}\right)$ as a constant value and simply multiply it to the sum

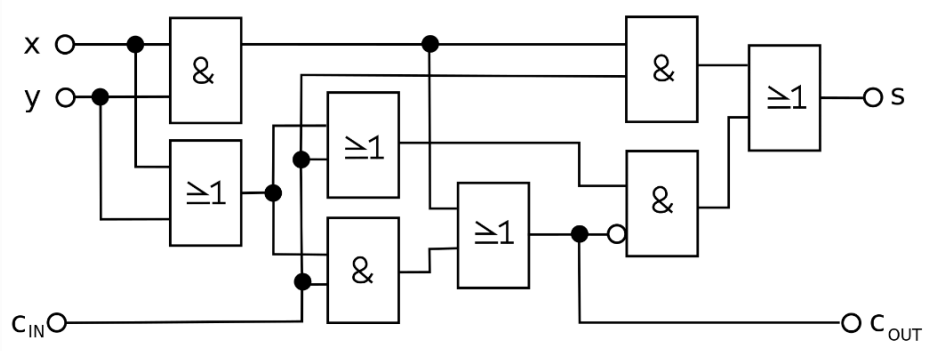
$$B\left(\frac{P_1}{P_1 + n_1}\right) = B\left(\frac{P_2}{P_2 + n_2}\right) = \dots = B\left(\frac{P_k}{P_k + n_k}\right)$$

$\sum P_k + n_k = P + n$ so this means

$$\sum \frac{P_k + n_k}{P + n} = \frac{P + n}{P + n} = 1$$

3.5

A full adder can be built using 4 AND gates, 4 OR gates, 1 NOT gate

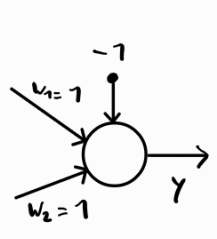


Perceptrons can be used to model the following logic gates: AND, OR, NOT

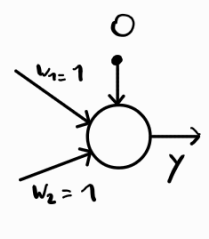
But a perceptron cannot express a XOR.

$$Activation\ function\ g(x) = \begin{cases} 1 & \text{if } x \geq 0.5 \\ 0 & \text{otherwise} \end{cases}$$

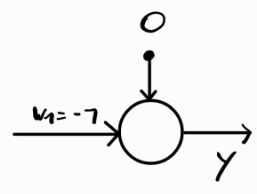
We define the logic gates as follows:



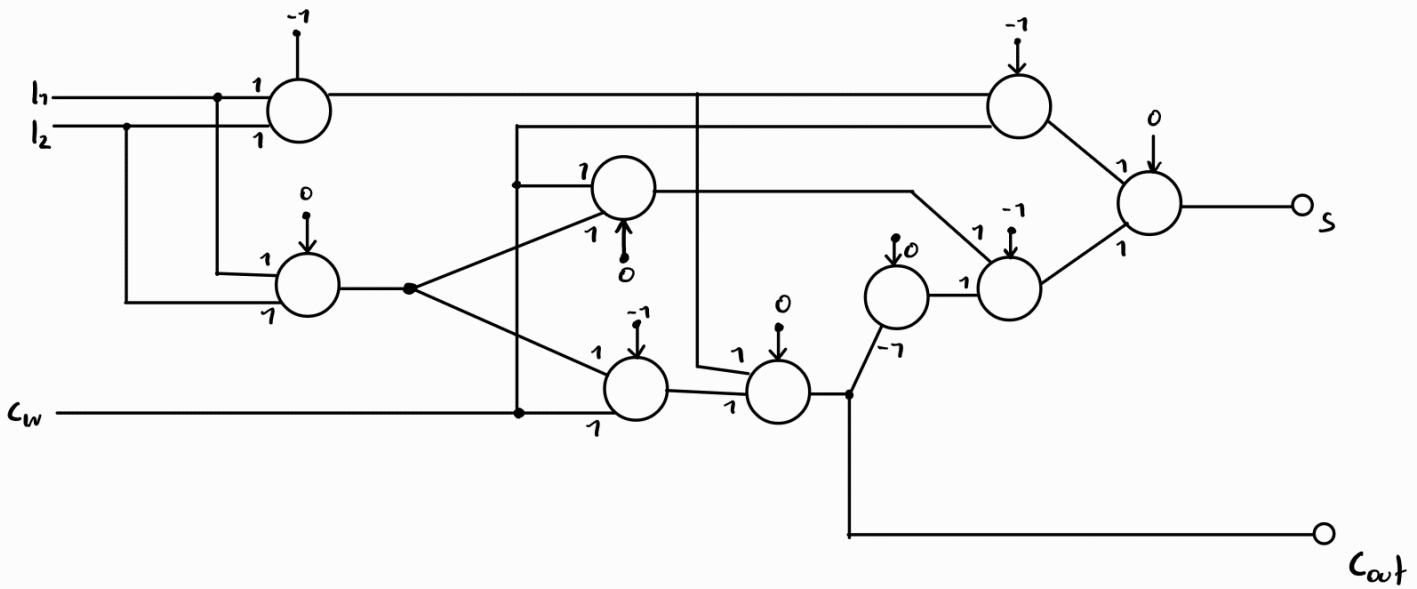
AND



OR



NOT



3.6

Variables/Definitions

Perceptron Learning Rule

$$w_i \leftarrow w_i + \alpha \cdot \underbrace{\text{Err} \cdot g'(in)}_{=\Delta} \cdot x_i$$

α ... learning rate (step size), parameter to control the weight search

$in = \sum_i w_i \cdot x_i = x \cdot w$ w ... weight vector
 x ... input vector
 y ... output

$h_w(x) = g(in)$

goal weights = -1, 7

$Err = y - h_w(x)$

Simplifications

$\alpha = 1, g(x) = x, g'(x) = 1$

$h_w(x) = g(in) = in$

$Err = y - in$

$w_i \leftarrow w_i + \alpha \cdot Err \cdot g'(in) \cdot x_i$

$w_i \leftarrow w_i + Err \cdot x_i$

① $f(1,2) = 2$:

$x_1 = 1 \quad w_1 = 0$

$x_2 = 2 \quad w_2 = 0$

$y = 2$

$in = 1 \cdot 0 + 2 \cdot 0 + 1$
bias

$Err = 2 - 1 = 1$

$w_1 \leftarrow 0 + 1 \cdot 1 = 1$

$w_2 \leftarrow 0 + 1 \cdot 2 = 2$

② $f(2,1) = 7$

$x_1 = 2 \quad w_1 = 1$

$x_2 = 1 \quad w_2 = 2$

$y = 7$

$in = 2 + 2 + 1 = 5$

$Err = 7 - 5 = 2$

$w_1 \leftarrow 1 + 2 \cdot 2 = 5$

$w_2 \leftarrow 2 + 2 \cdot 1 = 4$

③ $f(2,-1) = 4$

$x_1 = 2 \quad w_1 = 5$

$x_2 = -1 \quad w_2 = 4$

$y = 4$

$in = 10 - 4 + 1 = 7$

$Err = 4 - 7 = -3$

$w_1 \leftarrow 5 - 3 \cdot 2 = -1$

$w_2 \leftarrow 4 - 3 \cdot -1 = 7$

④ $f(8,2) = 7$

$x_1 = 8 \quad w_1 = -1$

$x_2 = 2 \quad w_2 = 7$

$y = 7$

$in = -8 + 14 + 1 = 7$

$Err = 7 - 7 = 0$

$w_1 \leftarrow -1 + 0 \cdot 8 = -1$

$w_2 \leftarrow 7 + 0 \cdot 2 = 7$

Final weights

$w_1 = -1$

$w_2 = 7$

bias = 1

$f(x_1, x_2) = -x_1 + 7x_2 + 1$

Check if learned weights are correct

$$f(x_1, x_2) = -x_1 + 7x_2 + 7$$

Expected: $f(1, 2) = 2$

Actual: $f(1, 2) = -1 + 7 \cdot 2 + 7 = 14 \quad \times$

Find solution via system of equations

I $1 \cdot w_1 + 2 \cdot w_2 + 7 = 2$

I $1 \cdot w_1 + 2 \cdot w_2 = 7$

II $2 \cdot w_1 + 1 \cdot w_2 + 7 = 7$

II $2 \cdot w_1 + 1 \cdot w_2 = 6$

III $2 \cdot w_1 - 1 \cdot w_2 + 7 = 4$

III $2 \cdot w_1 - 1 \cdot w_2 = 3$

IV $8 \cdot w_1 + 2 \cdot w_2 + 7 = 7$

IV $8 \cdot w_1 + 2 \cdot w_2 = 6$

I: $w_1 + 2w_2 = 7$

I in II:

$$w_1 = 7 - 2w_2$$

$$2(7 - 2w_2) + w_2 = 6$$

$$2 - 4w_2 + w_2 = 6$$

$$5w_2 = 4$$

$$w_2 = \frac{4}{5}$$

$$w_1 = 7 - 2 \cdot \frac{4}{5} = -\frac{3}{5}$$

w_1/w_2 in IV to verify:

$$8 \cdot -\frac{3}{5} + 2 \cdot \frac{4}{5} = 6$$

$$-\frac{24}{5} + \frac{8}{5} = 6$$

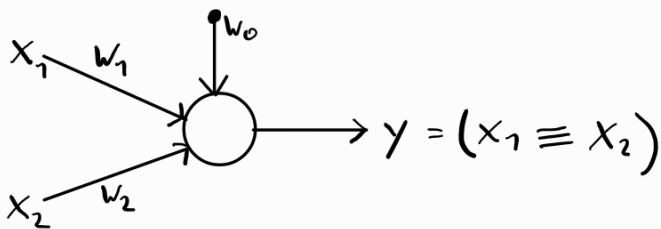
$$\frac{-16}{5} = 6 \quad \text{f.A } \times$$

no perfect weights can be found!

3.7

$$g(x) = \begin{cases} 1 & x \geq t \\ 0 & \text{otherwise} \end{cases}$$

t... some threshold



equivalence operator:

A	B	\equiv
0	0	1
0	1	0
1	0	0
1	1	1

A single-layer perceptron can be expressed in a single function:

$$f(x) = g\left(\sum_i x_i \cdot w_i\right) \quad f(x_1, x_2) = g(w_0 + x_1 \cdot w_1 + x_2 \cdot w_2)$$

According to the " \equiv " truth table the function expects the following input/output for the perceptron function:

$$f(0, 0) = 1$$

$$f(0, 1) = 0$$

$$f(1, 0) = 0$$

$$f(1, 1) = 1$$

We can represent this as a system of linear equations:

$$\begin{array}{l} \text{I } w_0 + 0 \cdot w_1 + 0 \cdot w_2 = 1 \\ \text{II } w_0 + 0 \cdot w_1 + 1 \cdot w_2 = 0 \\ \text{III } w_0 + 1 \cdot w_1 + 0 \cdot w_2 = 0 \\ \text{IV } w_0 + 1 \cdot w_1 + 1 \cdot w_2 = 1 \end{array} \longrightarrow \begin{array}{l} \text{I } w_0 = 1 \\ \text{II } w_0 + w_2 = 0 \\ \text{III } w_0 + 1 \cdot w_1 = 0 \\ \text{IV } w_0 + w_1 + w_2 = 1 \end{array} \longrightarrow \begin{array}{l} \text{II } 1 + w_2 = 0 \\ \text{III } 1 + w_1 = 0 \\ \text{IV } 1 + w_1 + w_2 = 1 \end{array}$$

$w_2 = -1, w_1 = -1$ in IV:

$$1 - 1 - 1 = 1$$

$$-1 = 1 \text{ f.A } \times$$

The system is incorrect, meaning no set of weights, that represent a logical equivalence operator can be found.

3.8

(1)

(2)

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

Sigmoid function σ is a softmax with $d=2$

-> cai-08-deep learning page 9/45

$$\text{softmax}(v)_k = \frac{e^{v_k}}{\sum_j e^{v_j}}$$

d ... size of vector v

v ... vector of values, $v = v_1, v_2, \dots, v_d$

examples:

$$v = (1, 2, 3) \quad \text{softmax}(v) = (0,001; 0,002; 0,997)$$

$$v = (5, 2, 0, -2) \quad \text{softmax}(v) = (0,946; 0,047; 0,006; 0,001)$$

if $d=2$ the two softmax values can be calculated as such:

$$\text{softmax}(v_1) = \frac{e^{v_1}}{e^{v_1} + e^{v_2}} \quad \text{softmax}(v_2) = \frac{e^{v_2}}{e^{v_1} + e^{v_2}}$$

$$\text{softmax}(v_1) = \frac{e^{v_1}}{e^{v_1} + e^{v_2}}$$

$$\text{softmax}(v_1) = \frac{e^{v_1}}{e^{v_1} \left(1 + \frac{e^{v_2}}{e^{v_1}}\right)}$$

$$\text{softmax}(v_1) = \frac{1}{1 + \frac{e^{v_2}}{e^{v_1}}}$$

$$\text{softmax}(v_1) = \frac{1}{1 + e^{v_2 - v_1}} \quad X = v_2 - v_1$$

$$\text{softmax}(\Delta v) = \frac{1}{1 + e^{-x}}$$