

Introduction to neural networks

Andreas Kleppe

Modification of slides made by Alexander Binder and Ole-Johan Skrede

27.01.2021



UiO : **Department of Informatics**
University of Oslo

- ⊙ For the time being, it seems that teachers can only participate on Padlet and Mattermost.
- ⊙ We encourage collaboration between students.
 - Explaining or discussing with fellow students are often more educational for all than to provide working code.
 - Students are allowed to answer each others requests on Padlet (or Mattermost).
 - The mandatory exercises and the exam must be individual work.
- ⊙ We aim to check Padlet more frequently.
- ⊙ Highly recommend doing the weekly exercises.

- 1 Classification by logistic regression
- 2 Artificial neurons
- 3 Neural networks
- 4 Example: The XOR problem
- 5 More on representability

Goal of classification: For every input sample $x \in \mathcal{X}$, correctly predict which class y it belongs to.

- ⊙ For 2-class classification, $y \in \{-1, +1\}$ or $y \in \{0, 1\}$.

Initial idea (from last lecture): Apply a linear mapping and classify according to the sign of the output:

$$f(x) = w \cdot x + b, \quad s(x) = \text{sign}(f(x))$$

While $f(x)$ has unbounded values, $s(x)$ is either -1 or 1 , but:

- ⊙ if $f(x) \approx 0$, we should be uncertain about the prediction.
- ⊙ if $f(x) \gg 0$, we should be confident about the prediction 1 .
- ⊙ if $f(x) \ll 0$, we should be confident about the prediction -1 .

How can we encode the uncertainty into a mapping from $(-\infty, +\infty)$ onto $[0, 1]$? We would like to map:

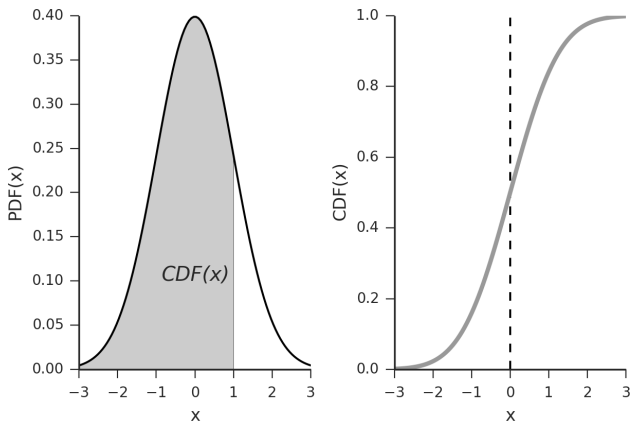
- ⊙ $f(x) \approx 0$ to $s(x) \approx 0.5$.
- ⊙ $f(x) \gg 0$ to $s(x) \approx 1$, in particular let $\lim_{x \rightarrow \infty} f(x) = 1$
- ⊙ $f(x) \ll 0$ to $s(x) \approx 0$, in particular let $\lim_{x \rightarrow -\infty} f(x) = 0$

This would allow us to interpret $s(x)$ as a probability.

There are many possible choices for the function $s(x)$.

We could e.g. use the cumulative distribution function (CDF) of any probability distribution function (PDF) with a median of 0.

- ⊙ The normal (Gaussian) distribution is used in probit (from **probability and unit**) regression.



We will use a simpler function called the logistic sigmoid function:

Definition: Logistic sigmoid function

$$s(x) = \frac{\exp(x)}{1 + \exp(x)} = \frac{1}{\exp(-x) + 1} \frac{\exp(x)}{\exp(x)} = \frac{1}{\exp(-x) + 1}$$

Note that:

$$s(0) = \frac{1}{\exp(-0) + 1} = \frac{1}{1 + 1} = 0.5$$

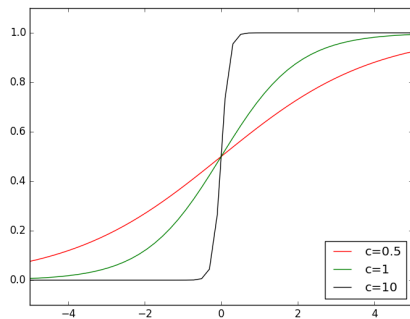
$$\lim_{x \rightarrow \infty} s(x) = \lim_{x \rightarrow \infty} \frac{1}{\exp(-x) + 1} = \frac{1}{0 + 1} = 1$$

$$\lim_{x \rightarrow -\infty} s(x) = \lim_{x \rightarrow -\infty} \frac{\exp(x)}{1 + \exp(x)} = \frac{0}{1 + 0} = 0$$

How does the logistic sigmoid function look?

Let's plot it for different scalings c :

$$s(cx) = \frac{\exp(cx)}{1 + \exp(cx)} = \frac{1}{\exp(-cx) + 1}$$



Definition: Logistic regression model

Assume we have a linear (or affine) mapping $f_{w,b}(x) = w \cdot x + b$. Plugging it into the logistic sigmoid function $s(x)$ provides a logistic regression model:

$$s(f_{w,b}(x)) = \frac{\exp(w \cdot x + b)}{1 + \exp(w \cdot x + b)} = \frac{1}{\exp(-w \cdot x - b) + 1}$$

For 2-class classification problems, the output $s(f_{w,b}(x)) \in [0, 1]$ is the predicted probability that sample x has class label $y = 1$, that is, $P(Y = 1|X = x)$.

A logistic model may also be called a logit (from **log**istic and **unit**) model.

Denoting the probability as p , the logit function is the logarithm of the odds:

$$\text{logit}(p) = \log \frac{p}{1-p}$$

The logit function is the inverse of the logistic sigmoid function:

$$s(\text{logit}(p)) = \frac{\exp(\text{logit}(p))}{1 + \exp(\text{logit}(p))} = \frac{\frac{p}{1-p}}{1 + \frac{p}{1-p}} = \frac{p}{1-p+p} = \frac{p}{1} = p$$

In statistics, logistic regression is often used to estimate the odds ratio between two events.

- ◉ We interpreted $s(f_{w,b}(x))$ as the probability of $Y = 1|X = x$ predicted by our model with parameters w and b .
- ◉ Given training samples x_1, \dots, x_n and associated class labels $y_1, \dots, y_n \in \{0, 1\}$, we would like to find suitable values for the model parameters w and b .
- ◉ If our logistic regression model was true, then the joint probability of observing the class labels would be:

$$P(Y_1 = y_1, \dots, Y_n = y_n | X_1 = x_1, \dots, X_n = x_n, w, b)$$

- ◉ Assuming that $(x_1, y_1), \dots, (x_n, y_n)$ are realisations of independent and identically distributed random variables $(X_1, Y_1), \dots, (X_n, Y_n)$, this reduce to:

$$\prod_{i=1}^n P(y_i | x_i, w, b)$$

- ◉ In maximum likelihood estimation, we consider this as a function of parameters w and b and maximise it.

The maximum likelihood estimates of the parameters w and b given the data points $(x_1, y_1), \dots, (x_n, y_n)$ are therefore:

$$(w^*, b^*) = \operatorname{argmax}_{(w,b)} \prod_{i=1}^n P(y_i | x_i, w, b)$$

- ⊙ If $y_i = 1$, then $P(y_i | x_i, w, b)$ is our model output $s(f_{w,b}(x_i))$.
- ⊙ If $y_i = 0$, then $P(y_i | x_i, w, b)$ is $1 - s(f_{w,b}(x_i))$.
- ⊙ This can be written as:

$$P(y_i | x_i, w, b) = s(f_{w,b}(x_i))^{y_i} (1 - s(f_{w,b}(x_i)))^{1-y_i}$$

Inserting this into the formula for the maximum likelihood estimates gives:

$$(w^*, b^*) = \operatorname{argmax}_{(w,b)} \prod_{i=1}^n s(f_{w,b}(x_i))^{y_i} (1 - s(f_{w,b}(x_i)))^{1-y_i}$$

Taking the logarithm does not change the location of the maximum (and it will make optimisation easier):

$$\begin{aligned}(w^*, b^*) &= \operatorname{argmax}_{(w,b)} \sum_{i=1}^n \log(s(f_{w,b}(x_i))^{y_i} (1 - s(f_{w,b}(x_i)))^{1-y_i}) \\ &= \operatorname{argmax}_{(w,b)} \sum_{i=1}^n y_i \log(s(f_{w,b}(x_i))) + (1 - y_i) \log(1 - s(f_{w,b}(x_i)))\end{aligned}$$

For convention, we will divide by n and minimise the negative:

$$(w^*, b^*) = \operatorname{argmin}_{(w,b)} \frac{1}{n} \sum_{i=1}^n -y_i \log(s(f_{w,b}(x_i))) - (1 - y_i) \log(1 - s(f_{w,b}(x_i)))$$

Thus, minimising this expression provides the maximum likelihood estimates of the weight parameter w and the bias parameter b in the logistic regression model $s(f_{w,b}(x))$ based on our data points $(x_1, y_1), \dots, (x_n, y_n)$.

This expression is called the cross-entropy loss of our data points and is the average cross-entropy loss of the individual data points.

Definition: Cross-entropy loss for 2-class classification

The cross-entropy loss of a data point (x_i, y_i) is defined as:

$$L(x_i, y_i | w, b) = -y_i \log(s(f_{w,b}(x_i))) - (1 - y_i) \log(1 - s(f_{w,b}(x_i)))$$

For one-hot labels $y_1, \dots, y_n \in \{0, 1\}$, the cross-entropy loss of a data point (x_i, y_i) is the negative logarithm of the predicted probability of the ground-truth class:

$$L(x_i, y_i | w, b) = \begin{cases} -\log(s(f_{w,b}(x_i))) & \text{if } y_i = 1 \\ -\log(1 - s(f_{w,b}(x_i))) & \text{if } y_i = 0 \end{cases}$$

Does such a log-probability make sense as a loss?

We derived that we want to minimise the loss function:

$$\begin{aligned}(w^*, b^*) &= \operatorname{argmin}_{(w,b)} \frac{1}{n} \sum_{i=1}^n -y_i \log(s(f_{w,b}(x_i))) - (1 - y_i) \log(1 - s(f_{w,b}(x_i))) \\ &= \operatorname{argmin}_{(w,b)} \frac{1}{n} \sum_{i=1}^n L(x_i, y_i | w, b)\end{aligned}$$

Let's use gradient descent.

The gradient of the logistic sigmoid function is:

$$\begin{aligned}\frac{\partial s(x)}{\partial x} &= \frac{\partial}{\partial x} \frac{1}{\exp(-x) + 1} = \frac{\exp(-x)}{(\exp(-x) + 1)^2} \\ &= \frac{1}{\exp(-x) + 1} \left(\frac{\exp(-x) + 1 - 1}{\exp(-x) + 1} \right) = s(x)(1 - s(x))\end{aligned}$$

Thus:

$$\begin{aligned}\frac{\partial \log(s(x))}{\partial x} &= \frac{1}{s(x)} \frac{\partial s(x)}{\partial x} = 1 - s(x) \\ \frac{\partial \log(1 - s(x))}{\partial x} &= -\frac{1}{1 - s(x)} \frac{\partial s(x)}{\partial x} = -s(x)\end{aligned}$$

With $f_{w,b}(x) = w \cdot x + b$, we get:

$$\begin{aligned}\frac{\partial L}{\partial w} &= \frac{\partial}{\partial w} \frac{1}{n} \sum_{i=1}^n -y_i \log(s(f_{w,b}(x_i))) - (1 - y_i) \log(1 - s(f_{w,b}(x_i))) \\ &= \frac{\partial}{\partial w} \frac{1}{n} \sum_{i=1}^n -y_i \log(s(w \cdot x_i + b)) - (1 - y_i) \log(1 - s(w \cdot x_i + b)) \\ &= \frac{1}{n} \sum_{i=1}^n -x_i y_i (1 - s(w \cdot x_i + b)) + x_i (1 - y_i) s(w \cdot x_i + b) \\ &= \frac{1}{n} \sum_{i=1}^n -x_i y_i + x_i s(w \cdot x_i + b) = \frac{1}{n} \sum_{i=1}^n x_i (s(w \cdot x_i + b) - y_i) \\ \frac{\partial L}{\partial b} &= \frac{1}{n} \sum_{i=1}^n s(w \cdot x_i + b) - y_i\end{aligned}$$

The gradient of a data point (x_i, y_i) is thus given by:

$$\frac{\partial L(x_i, y_i | w, b)}{\partial w} = x_i (s(w \cdot x_i + b) - y_i)$$
$$\frac{\partial L(x_i, y_i | w, b)}{\partial b} = s(w \cdot x_i + b) - y_i$$

Observation 1:

Nice interpretation: The partial derivatives are the difference between the predicted value $s(w \cdot x_i + b)$ and the true value $y_i \in \{0, 1\}$, with or without weighting by the sample x_i .

Observation 2:

If both classes in the data can be perfectly separated by a linear mapping, then optimisation will try to make the weights w to go to infinity. This is a source of instability during optimisation time.

- ⊙ If $y_i = 1$, then the loss for the data point (x_i, y_i) is $-\log s(f_{w,b}(x_i))$.
- ⊙ For $-\log s(f_{w,b}(x_i))$ to be 0, then $s(f_{w,b}(x_i))$ need to be 1.
- ⊙ For $s(f_{w,b}(x_i)) = \frac{1}{\exp(-f_{w,b}(x_i))+1}$ to be 1, then $\exp(-f_{w,b}(x_i))$ need to be 0.
- ⊙ For $\exp(-f_{w,b}(x_i))$ to be 0, then $f_{w,b}(x_i) = w \cdot x_i + b \rightarrow \infty$.
- ⊙ For data points where $y_i = 0$, a loss $-\log(1 - s(f_{w,b}(x_i))) = 0$ implies that $w \cdot x_i + b$ needs to go to $-\infty$.
- ⊙ For linearly separable problems, the total loss will thus approach 0 as $|w|$ approach infinity (and is appropriately oriented). The optimisation may therefore result in numerical instability.
- ⊙ So making $P(Y = y_i | X = x_i) = 1$ for all data points requires increasing $|w|$ to a very large value.

Why this does not happen when the training dataset is not linearly separable?

- ⊙ No matter the orientation of w , there will be always a data point (x_i, y_i) that is misclassified.
- ⊙ Suppose $y_i = 1$ so that the loss of this sample is $-\log s(f_{w,b}(x_i))$.
- ⊙ We can simulate $|w| \rightarrow \infty$ by using as weight cw and let go $c \rightarrow \infty$.
- ⊙ Because (x_i, y_i) is misclassified, $f_{w,b}(x_i) = cw x_i + b < 0$, and thus $\lim_{c \rightarrow \infty} cw x_i + b = -\infty$.
- ⊙ Therefore, the predicted probability of the correct class $s(f_{w,b}(x_i))$ will go to 0, making the loss go to infinity.
- ⊙ Because the loss is always non-negative, this single misclassification would make the total loss go to infinity if $|w|$ goes to infinity.
- ⊙ Thus, the absolute value of the weight will not go to infinity if the training dataset is not linearly separable.
- ⊙ There will instead be a trade-off between pushing $P(Y = y_i | X = x_i)$ close to 1 for corrected classified samples and not pushing it close to 0 for misclassified samples.

Key takeaways:

- Applying the logistic sigmoid function to a linear mapping provides a logistic regression model.
- Outputs a score in $[0, 1]$ which can be interpreted as the predicted probability of the class labelled 1.
- One minus score is our predicted probability of class label 0.
- Logistic regression is therefore a classification algorithm.
- The cross-entropy loss can be derived as the function to optimise by maximum likelihood estimation.
- The cross-entropy loss is the sum of the negative logarithm of the predicted probabilities weighted by the ground-truth distribution.
- For one-hot labels, this becomes the negative logarithm of the predicted probability of the ground-truth class.

- ① Classification by logistic regression
- ② Artificial neurons**
- ③ Neural networks
- ④ Example: The XOR problem
- ⑤ More on representability

A mathematical function that for instance can:

- ⊙ Take some inputs $\{z_i\}_{i=1}^d$.
- ⊙ Depend on some weights w_{ij} and bias b_j .
- ⊙ Apply some non-linear activation function $g(\cdot)$.
- ⊙ Output z_j .

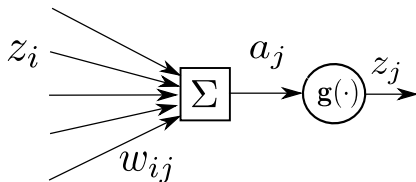
Forward equation:

$$a_j = \sum_{i=1}^d z_i w_{ij} + b_j$$

linear mapping

$$z_j = g(a_j)$$

nonlinear activation

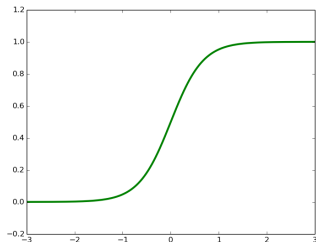


- ⊙ A function $g : \mathbb{R}^1 \rightarrow \mathbb{R}^1$
- ⊙ Introduces non-linearity.
- ⊙ Should be differentiable if you are planning to use gradient-based optimisation.
- ⊙ Typically prevents or reduces weak signals from passing through.
- ⊙ Intuition: Determines what the neuron fires for different outputs of the linear mapping.

We could for instance use the logistic sigmoid function:

$$g(x) = \frac{1}{\exp(-x) + 1}$$

- ⦿ Called sigmoid activation.
- ⦿ Historically popular.
- ⦿ Currently rarely used.



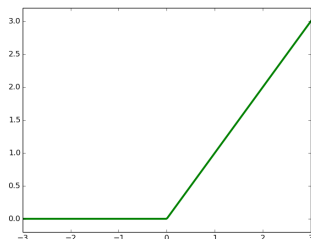
Suffer from gradient saturation: its derivative gets exponentially small for large $|x|$:

$$\frac{\partial s(x)}{\partial x} = \frac{e^{-x}}{(e^{-x} + 1)^2} \leq \begin{cases} e^{-x} & x \geq 0 \\ \frac{1}{e^{-x} + 1} & x < 0 \end{cases}$$

An alternative is to simply set negative values to 0, i.e. let:

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

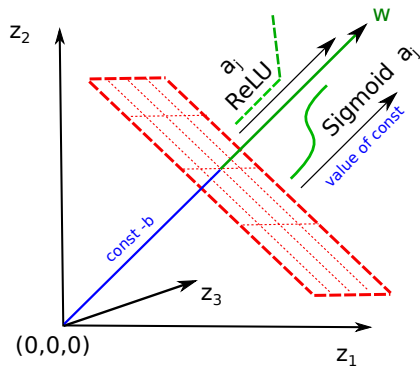
- ⊙ Called ReLU activation.
- ⊙ Gradient do not saturate.
- ⊙ Not differentiable in 0.
 - Not of practical concern.
 - Could define it to be 1.



The gradient is 0 for negative inputs.

- ⊙ Might cause inactive nodes to remain inactive.
- ⊙ Might still work fine in practice.
- ⊙ Careful selection of learning rate may be important.
- ⊙ There exists alternatives with non-zero gradient for negative inputs, e.g. leaky ReLU: $g(x) = \max(0.01x, x)$

A popular choice, particularly with convolutional neural networks.

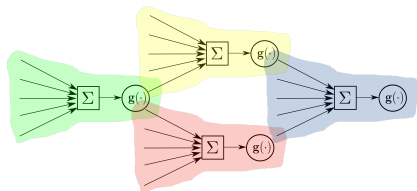


- ⊙ The weights w defines the orientation of the plan.
- ⊙ The bias b defines the position of the plan.
- ⊙ Output of activation function is constant in the red plane.
- ⊙ Output of activation function varies along the direction of w .

- ① Classification by logistic regression
- ② Artificial neurons
- ③ Neural networks**
- ④ Example: The XOR problem
- ⑤ More on representability

A neural network is a directed graph structure made from connected neurons.

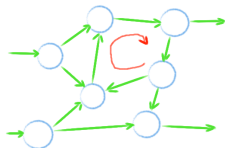
- A neuron can be input to many other neurons.
- A neuron can receive input to many other neurons.
- Each neuron has same structure ($g(\cdot), \Sigma$) but different parameters (w_{ij}, b_j).
- Can stack neurons in layers.



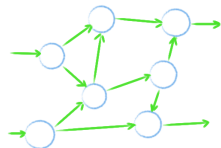
Definition: Neural Network

Any directed graph built from neurons is a neural network.

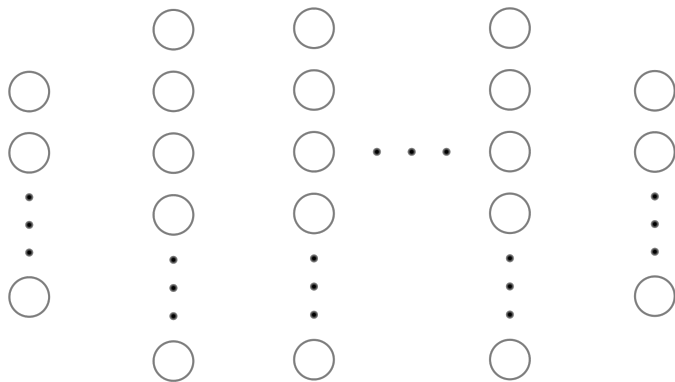
Two important types: recurrent and feedforward neural networks

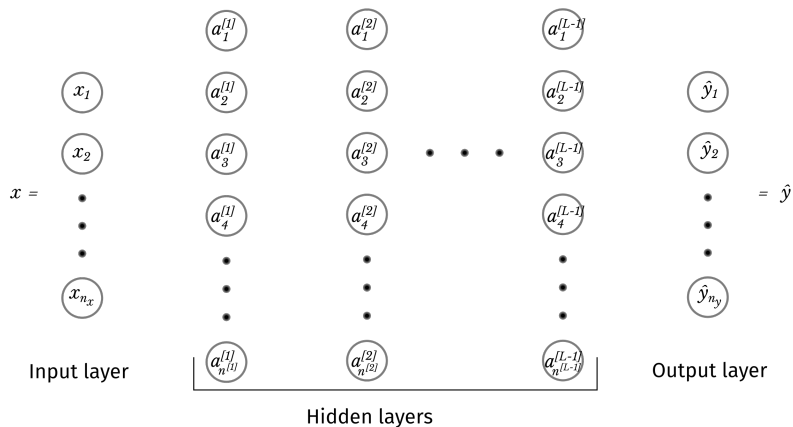


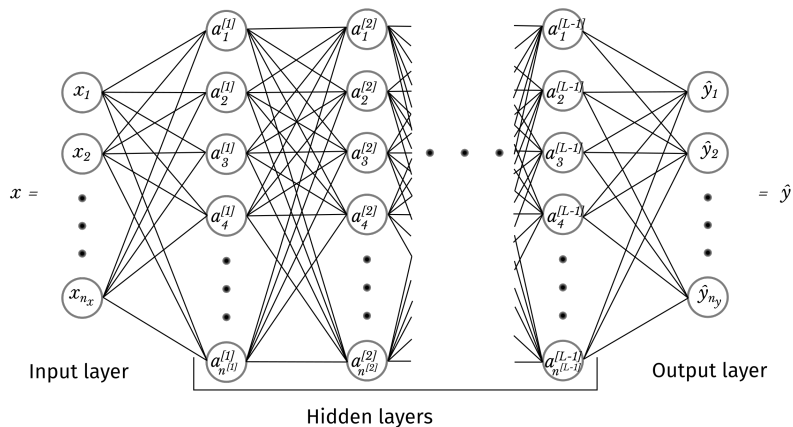
recurrent (not covered in this lecture)
e.g. for speech processing



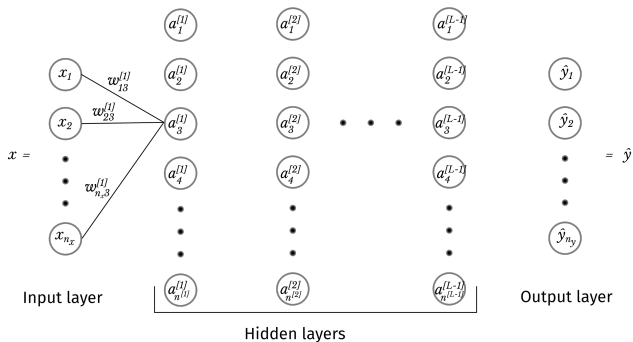
feedforward
e.g. for image classification



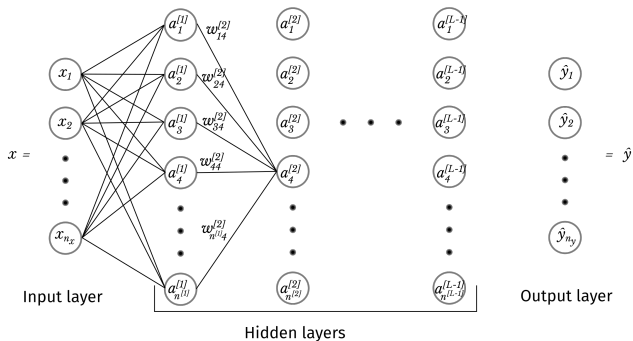




- Superscript with square brackets $[l]$: layer l
 - L : Number of layers in the network.
 - $n^{[l]}$: Number of nodes in layer l
 - $n_x = n^{[0]}$: Input dimension
 - $n_y = n^{[L]}$: Output dimension (number of classes)
 - x : Array of inputs
 - y : Array of *true* outputs
 - \hat{y} : Array of *predicted* output
- w : Edge weights
 - b : Node bias
 - z : Linear combination of activations from previous layer
 - $a^{[l]}$: Node activation in layer l .
 - $a^{[0]} = x$: Input vector
 - $a^{[L]} = \hat{y}$: Output vector
 - Subscript j or jk : Element in vector or matrix
 - m : Number of samples in the training dataset



$$a_3^{[1]} = g \left(\sum_{j=1}^{n_x} w_{j3}^{[1]} x_j + b_3^{[1]} \right)$$

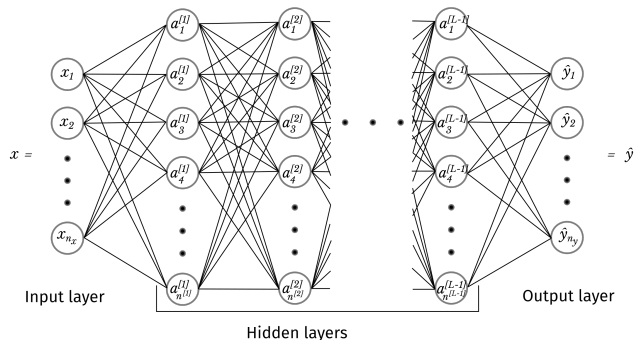


$$a_4^{[2]} = g \left(\sum_{j=1}^{n^{[1]}} w_{j4}^{[2]} a_j^{[1]} + b_4^{[2]} \right)$$

- $a_k^{[l]}$ is the *activation* of node k in layer l :

$$a_k^{[l]} = g \left(\sum_{j=1}^{n^{[l-1]}} w_{jk}^{[l]} a_j^{[l-1]} + b_k^{[l]} \right)$$

- $w_{jk}^{[l]}$ is the *weight* from node j in layer $l - 1$ to node k in layer l .
- $b_k^{[l]}$ is the *bias* of node k in layer l .
- All above are scalars.
- g is the activation function.
- All w and b are “trainable”, and will be adjusted according to some optimization routine.
- By convention:
 - $a_k^{[0]} = x_k$ and $a_k^{[L]} = \hat{y}_k$
 - The network have L layers (we do not include the input layer in the count) and $L - 1$ hidden layers.



$$a_k^{[l]} = g \left(\sum_{j=1}^{n^{[l-1]}} w_{jk}^{[l]} a_j^{[l-1]} + b_k^{[l]} \right), k \in \{1, 2, \dots, n^{[l]}\}, l \in \{1, 2, \dots, L\}$$

- ⊙ Applying an activation function at the last layer will in general not provide outputs that can be interpreted as probabilities.
- ⊙ For 2-class classification, we could have applied the logistic sigmoid function instead of an activation function to produce probabilities from the single output neuron.
- ⊙ We will instead apply a generalisation of the logistic sigmoid function called the softmax function:

$$s(x)_k = \frac{e^{x_k}}{\sum_{i=1}^n e^{x_i}}$$

- ⊙ This function converts the linearly mapped values in the output layer, z_1, \dots, z_{n_y} , to predicted probabilities.

In the output layer we have:

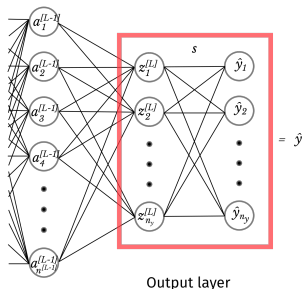
$$z_k^{[L]} = \sum_{j=1}^{n^{[L-1]}} w_{jk}^{[L]} a_j^{[L-1]} + b_k^{[L]}$$

$$\begin{aligned} a_k^{[L]} &= s(z_k^{[L]})_k \\ &= \hat{y}_k \end{aligned}$$

for:

$$\begin{aligned} k &= 1, \dots, n_y, \\ &= 1, \dots, n^{[L]}. \end{aligned}$$

$z^{[L]}$ are called *logits*, and \hat{y} will be interpreted as the predicted probabilities.



Numerical instability can be a problem when using the softmax function, because of the exponential function and division.

Two common “tricks” that can help:

1. Shift exponential arguments to max zero:

$$s(x)_k = \frac{e^{x_k}}{\sum_{i=1}^n e^{x_i}} = \frac{e^{x_k - \max(x)}}{\sum_{i=1}^n e^{x_i - \max(x)}}$$

2. Take logarithm and exponentiate it to get rid of division:

$$t(x)_k = \log s(x)_k = x_k - \log \sum_{i=1}^n e^{x_i}$$
$$s(x)_k = e^{t(x)_k}$$

It is possible to combine these two “tricks”.

The derivation and definition of cross-entropy loss can be extended to any model's predicted probabilities $\hat{y}_i \in [0, 1]^{n_y}$ associated with input sample $x_i \in \mathbb{R}^{n_x}$ with true class label $y_i \in [0, 1]^{n_y}$:

Definition: Cross-entropy loss

The cross-entropy loss of a data point (x_i, y_i) is defined as:

$$L(\hat{y}_i, y_i) = - \sum_{k=1}^{n_y} (y_i)_k \log(\hat{y}_i)_k$$

For one-hot labels, i.e. each y_i is 0 in all elements except one where it is 1, the cross-entropy loss of a data point (x_i, y_i) is the negative logarithm of the predicted probability of the ground-truth class:

$$L(\hat{y}_i, y_i) = - \log(\hat{y}_i)_k \text{ where } y_i = e_k$$

- As for the logistic regression model, we can use gradient descent to optimise the cross-entropy loss function:

$$\Theta^* = \operatorname{argmin}_{\Theta} \left\{ -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^{n_y} (y_i)_k \log(\hat{y}_i)_k \right\}.$$

- However, we have a lot more parameters this time:

$$\Theta = \{w_{jk}^{[l]}, b_k^{[l]} : j \in \{1, \dots, n^{[l-1]}\}, k \in \{1, \dots, n^{[l]}\}, l \in \{1, \dots, L\}\}$$

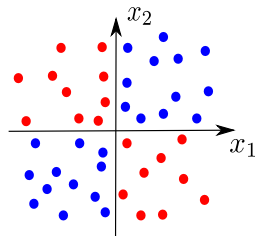
- How to do this in practice will be the subject of a later lecture.

Key takeaways:

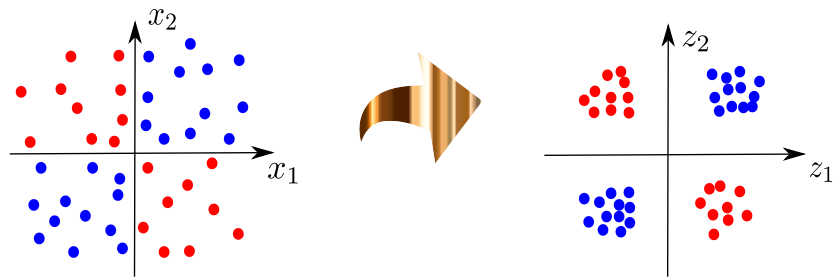
- ⦿ Logistic regression is an exemplary model for a single neuron network, that is, a 1-layer neural network (0 hidden layers) with one node in the output layer.
- ⦿ Stacking multiple layers gives a deep neural network.
- ⦿ Each neuron (in hidden layers) in dense feedforward neural networks receives input from all neurons in the previous layer and sends its activation (output) to all neurons in the next layer.
- ⦿ Note: The particular choice of network architecture is often subordinate to the selecting appropriate loss function, algorithm to updating weights, data augmentation and more.

- ① Classification by logistic regression
- ② Artificial neurons
- ③ Neural networks
- ④ Example: The XOR problem
- ⑤ More on representability

To illustrate how concatenation of neurons allows learning non-linear mappings, consider separating red from blue samples in the following example:



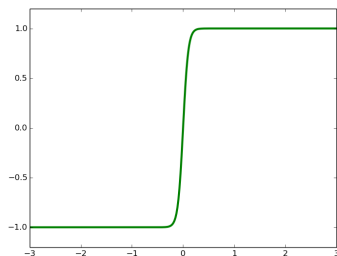
Samples lie in quadrants around coordinated $(\pm 1, \pm 1)$.

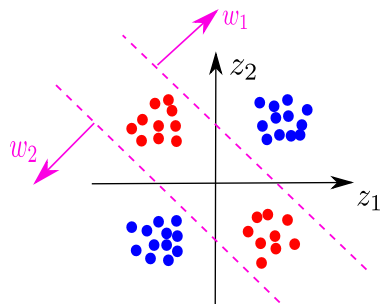


$$z_1 = \tanh(10x_1)$$

$$z_2 = \tanh(10x_2)$$

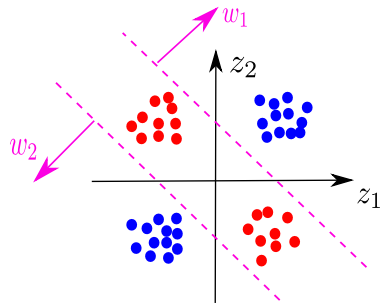
- ⊙ \tanh pushes points towards $(\pm 1, \pm 1)$



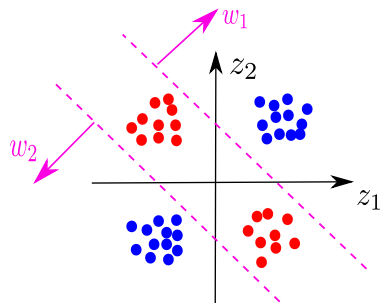


Idea:

- ⊙ For proper choice of w_1 and b_1 , $g_1 = \sigma(w_1 z + b_1)$ has high values ≈ 1 in the upper right $z \approx (1, 1)$ and low values ≈ 0 elsewhere.
- ⊙ For proper choice of w_2 and b_2 , $g_2 = \sigma(w_2 z + b_2)$ has high values ≈ 1 in the lower left $z \approx (-1, -1)$ and low values ≈ 0 elsewhere.
- ⊙ Then, $g_1 + g_2$ will be ≈ 1 in the upper right and lower left corners.
- ⊙ In the middle zone, $g_1 + g_2$ will be ≈ 0 .



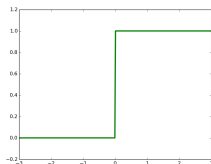
⊙ $w_1 = ?$, $w_2 = ?$



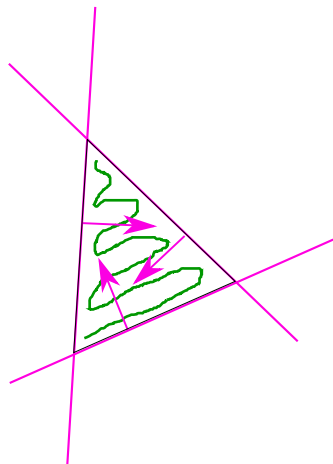
- ⊙ If $w_1 = (1, 1)$, $w_2 = (-1, -1)$, then $w_1 \cdot z \approx 0$, $w_2 \cdot z \approx 0$ for middle zone, while $w_1 \cdot z \approx 2$ in the upper right corner and $w_2 \cdot z \approx 2$ in the lower left corner.
- ⊙ Could then choose $-2 < b_1 < 0$ and $-2 < b_2 < 0$.
- ⊙ Final output: $f(x) = \sigma(c(g_1 + g_2) + b_3)$

- ① Classification by logistic regression
- ② Artificial neurons
- ③ Neural networks
- ④ Example: The XOR problem
- ⑤ More on representability

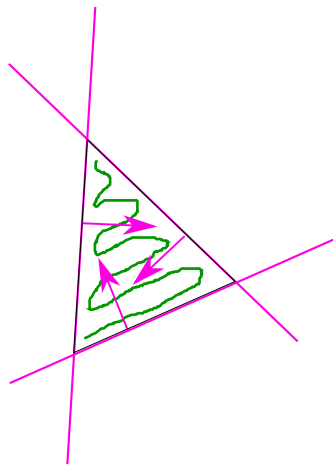
- Use threshold activation (for simplicity):



- A 2-layer network with n neurons in the hidden layer, and a thresholded sum in the second layer can represent any convex polygonal shape with n edges.
- With $n \rightarrow \infty$ neurons in the first layer, any convex shape can be **approximately** encoded.



- Adding a third layer allows the neural network to approximately encode any union of convex shapes.
- However, being able to (approximately) represent does not imply that this would be the result when learning from finite data.



Universal approximation theorem (a version)

Let $g(\cdot)$ be a continuous function on a m -dimensional hypercube $[0, 1]^m$. Let $a(\cdot)$ be a non-constant, bounded, continuous (activation) function. Then $g(\cdot)$ can be approximated arbitrarily well, that is, for every maximal deviation $\epsilon > 0$, there exists a set of weights u_i, w_i and biases b_i such that:

$$\forall x \in [0, 1]^m : |g(x) - (\sum_{i=1} u_i a(w_i x + b_i) + b)| < \epsilon$$

Neural networks **with one hidden layer and two layers of weights** can approximate any smooth function on a compact hypercube. If there is a good algorithm for learning the parameters from finite data, we are done!

- ⊙ Universal approximation theorem was the arch-seducer of neural network research.
- ⊙ Kolmogorov (1957), Hornik (1989), Cybenko (1989) and others: Neural network can approximate *any continuous function on a compact region*.
- ⊙ However, ability to approximate any function \neq able to learn well from finite training data

Key learning goals

- ⦿ Understand logistic regression and why it is a classification algorithm and a 1-layer neural network.
- ⦿ Know the structure and components of dense feedforward neural networks.
- ⦿ Be able to do forward pass by hand for dense feedforward neural networks.
- ⦿ Understand the cross-entropy loss.
- ⦿ Have an intuition about representability of simple neural networks.

Questions?