

SCHOOL OF ELECTRONIC ENGINEERING AND COMPUTER SCIENCE
QUEEN MARY UNIVERSITY OF LONDON

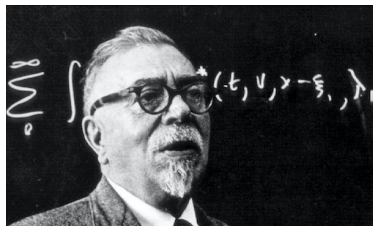
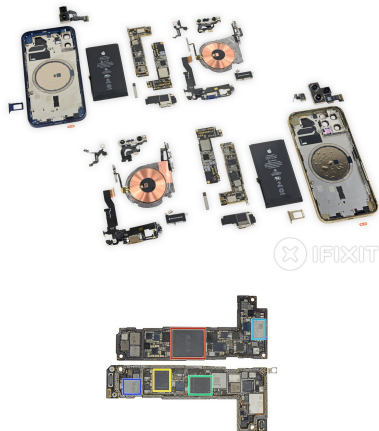
ECS7020P Principles of Machine Learning

Neural Networks and Deep Learning

Dr Jesús Requena Carrión

16 Nov 2021

Matched and Wiener filtering



Data never beats what is known!

Remember linear boundaries?

Linear classifiers use boundaries defined by the linear equation $\mathbf{w}^T \mathbf{x} = 0$.

In a binary problem with classes 0 and 1, any new sample \mathbf{x}_i can be classified as follows:

- If $\mathbf{w}^T \mathbf{x}_i > 0$, the sample belongs to class 1.
- If $\mathbf{w}^T \mathbf{x}_i < 0$, the sample belongs to class 0.

Non-linear boundaries

We can create non-linear boundaries using the same principle. A boundary is defined by an equation $f(\mathbf{x}) = 0$ and for a new \mathbf{x}_i we do the following:

- If $f(\mathbf{x}_i) > 0$, the sample belongs to class 1.
- If $f(\mathbf{x}_i) < 0$, the sample belongs to class 0.

For instance, QDA defines quadratic boundaries, that can be expressed as $f(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 + w_3x_1x_2 + w_4x_1^2 + w_5x_2^2$ in a 2D predictor space.

Using linear boundaries to create non-linear ones

Can we use linear boundaries to create a non-linear boundary?

Agenda

Neural networks

What can perceptrons do?

Neural networks as machine learning models

Types of layers

Summary

Neural networks: Not quite new

Neural networks were first proposed in the 1940s, became very popular in the 70s and 80s and fell out of favour in the 90s, as:

- Other algorithms performed better.
- They have a large number of parameters and are hard to train.
- Analysing them is difficult.

From 2010 onwards there has been a resurgence of interest in neural networks because of:

- Improved **optimisation** algorithms, increased **computational power** and **larger datasets**.
- Careful definition of **architectures**.
- Application of **transfer learning**.

What is a neural network?

A neural network is a **computing system** loosely inspired by the human nervous system, commonly used as a **family** of Machine Learning models. We will focus on their application to **supervised problems**.

From a **functional** point of view, a neural network:

- Produces an **output** (label) given an **input** (predictors), .
- Has **weights** (parameters) that can be **tuned** (trained).



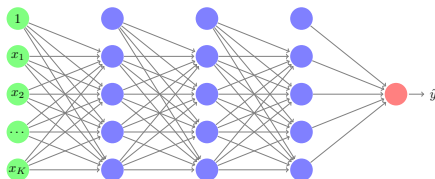
W denotes all the weights of the neural network and $h_W(x)$ indicates explicitly that the prediction \hat{y} depends on W .

What is a neural network?

From an **computational** angle, neural networks consist of interconnected units that (loosely) mimic **neurons**. This architecture is appealing since:

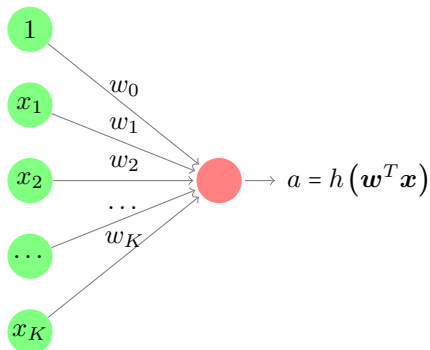
- **Neuroscience** suggests biological neural networks can solve any problem.
- **Mathematics** suggests artificial neural networks can reproduce any input/output relationship, provided they are complex enough.

Hence, the family of neural network models can be seen as a **universal machine**.



The perceptron

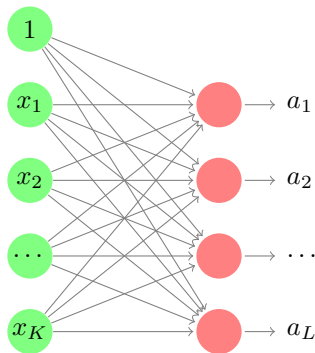
The perceptron is the **basic unit** of a neural network. It is defined by a **weight vector** w and an **activation function** $h(\cdot)$ that map an extended vector x to an output a . The coefficient w_0 is known as the **bias**.



The activation function is **non-linear**.

Layers

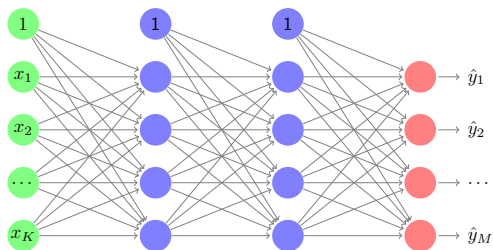
A **layer** is a collection of perceptrons that use the **same input**. Each perceptron within a layer produces a separate output.



A layer consisting of L perceptrons and K inputs has $L \times (K + 1)$ weights.

The architecture

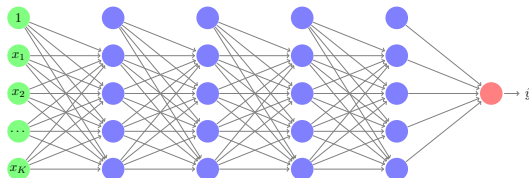
The **architecture** of a neural network describes how layers are connected. The **input layer** is the predictor vector, **hidden layers** produce internal features and the **output layer** produces the prediction.



Vector $\hat{\mathbf{y}} = [\hat{y}_1, \dots, \hat{y}_M]^T = h_{\mathbf{W}}(\mathbf{x})$ is the output of the neural network.
(Note that here \hat{y}_j is not the prediction for the j -th sample!)

The neural network

- A neural network consists of **perceptrons** arranged in **layers** that are connected according to an **architecture**.
- There is one **parameter** per connection (the connections's **weight**).
- Each layer produces intermediate **features**.
- The number of layers determines the depth of the neural network (hence the distinction between **shallow** and **deep** neural networks).



Agenda

Neural networks

What can perceptrons do?

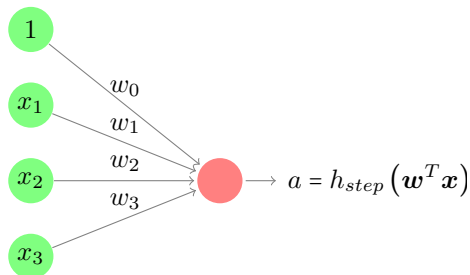
Neural networks as machine learning models

Types of layers

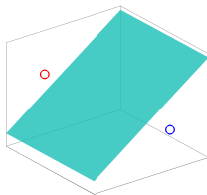
Summary

The perceptron as a linear classifier

A perceptron that uses a **step** activation function is actually a **linear classifier**. What if it uses a **logistic function**?



$$h_{step}(d) = \begin{cases} 1, & \text{if } d > 0. \\ 0, & \text{otherwise.} \end{cases}$$

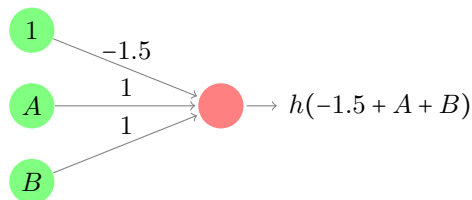


In a 3D predictor space, a perceptron with a step activation function separates two decision regions by a plane!

The perceptron as a basic logical function

We can use perceptrons to implement logical functions.

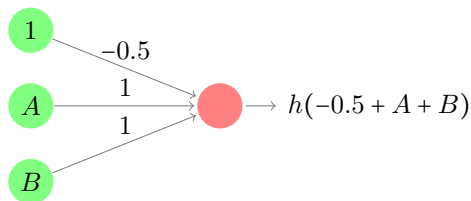
In the perceptron below, assuming that the activation function is a step function and A and B can be either 0 or 1, the output $f(A, B)$ corresponds to a logical AND of A and B .



| A | B | $f(A, B)$ |
|-----|-----|-----------|
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

The perceptron as a basic logical function

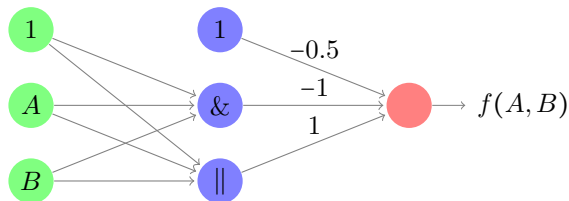
What logical function would you obtain if the coefficients were $w_0 = -0.5$, $w_A = 1$ and $w_B = 1$ instead?



| A | B | $f(A, B)$ |
|-----|-----|-----------|
| 0 | 0 | |
| 0 | 1 | |
| 1 | 0 | |
| 1 | 1 | |

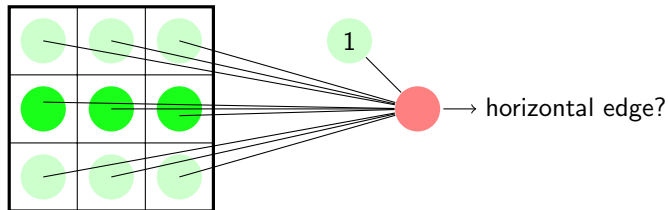
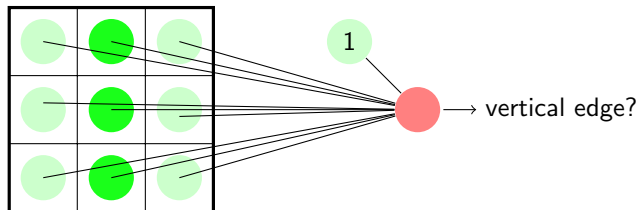
Adding depth to derive new logical functions

What about the following neural network? In this diagram, $\&$ denotes an AND perceptron, and \parallel an OR perceptron.

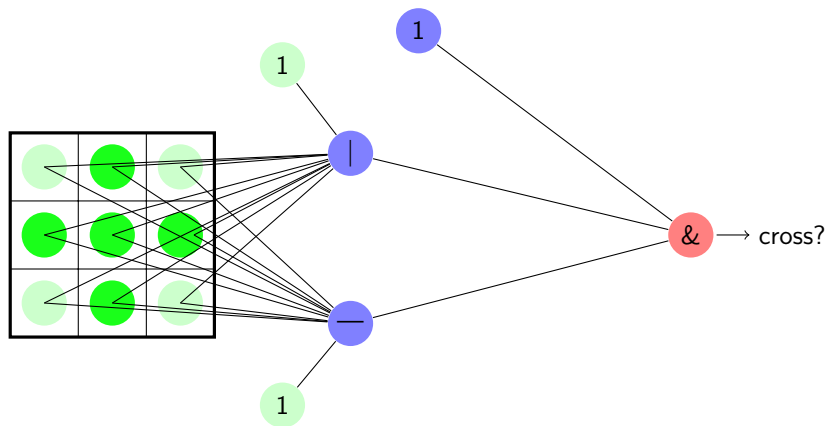


| A | B | $f(A, B)$ |
|-----|-----|-----------|
| 0 | 0 | |
| 0 | 1 | |
| 1 | 0 | |
| 1 | 1 | |

The perceptron as a grid pattern detector

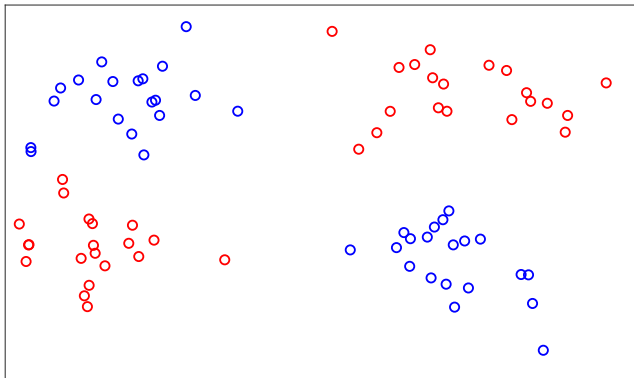


Adding depth to detect derived grid patterns



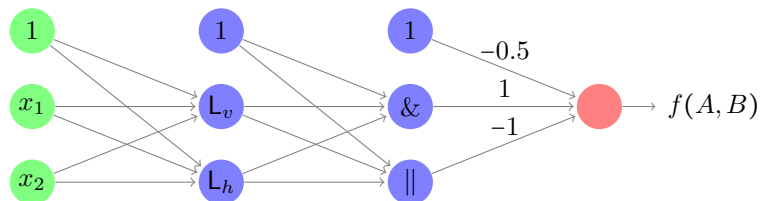
Combining linear classifiers and logical functions

Build a neural network classifier for the following dataset.



Combining linear classifiers and logical functions

In this diagram, L_v and L_h represent perceptrons that implement a vertical and horizontal linear boundaries, respectively. Their outputs indicate whether samples are on one side of the boundary or another.



Versatility of neural networks as a computing system

Using a single perceptron we can implement:

- Linear boundaries.
- Logical functions.
- Grid pattern detectors.

Connecting perceptrons allows us to implement more complex operations, e.g. complex boundaries, complex logical functions and complex grid patterns. This can be seen as building complexity from basic operations.

So far we have seen neural networks from a computation angle and have specified each network by ourselves. Note this is **not** machine learning!

Agenda

Neural networks

What can perceptrons do?

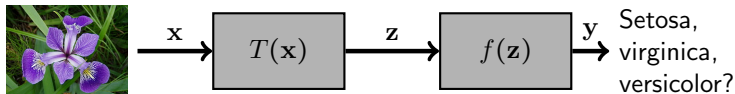
Neural networks as machine learning models

Types of layers

Summary

Machine learning pipelines: Reminder

A machine learning pipeline is a sequence of data operations.



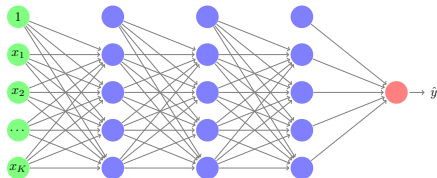
A simple pipeline consists of a first transformation stage followed by a machine learning model:

- The transformation stage produces **derived** features.
- The model uses the derived features as input.
- The transformation stage can also be trained.

Neural networks as tunable pipelines

A neural network can be seen as an entire tunable machine learning pipeline, where:

- Each perceptron defines one **derived feature** or **concept** using other features, raw or derived.
- **Increasing the number of perceptrons** in a layer allows to create more new concepts per layer.
- **Increasing the number of layers** (deeper networks) allows to create concepts of increasing complexity.



Deep neural networks: Computational angle

From a cognitive point of view, large neural networks are appealing, as they give us the necessary **flexibility to create new and increasingly complex concepts** that might be relevant to make a prediction.

However:

- Higher flexibility increases the risk of **overfitting** (which we can see as a network creating and using irrelevant concepts).
- Large number of parameters need to be tuned, therefore the **computational requirements** might be too high.

For **complex inputs**, such as pictures consisting of millions of pixels, this is even more severe.

Training neural networks: Cost function

Every Machine Learning algorithm needs a **model**, a **cost function** and an **optimisation** method.

Given a dataset $\{(\mathbf{x}_i, y_i), 1 \leq i \leq N\}$, where labels can take on the values 0 or 1, a common cost function for classification is the **negative log-likelihood function**, defined as:

$$l(\mathbf{W}) = -\frac{1}{N} \sum_{n=1}^N y_i \log [\hat{y}_i] + (1 - y_i) \log [1 - \hat{y}_i]$$

where $\hat{y}_i = h_{\mathbf{W}}(\mathbf{x}_i)$. This cost function can be extended to multi-class classifiers.

Gradient descent and back-propagation

Gradient descent is the method of choice to find the optimal set of coefficients \mathbf{W} for the cost function $l(\mathbf{W})$.

Obtaining the gradient is easy, but can be computationally expensive.

Back-propagation is an efficient algorithm to **compute the gradient**. This gradient is then used by the optimisation algorithm to update \mathbf{W} .

Back-propagation exploits the **chain rule of calculus**. It turns out that to compute the gradient in one layer, we just need information from the next layer. Back-propagation starts from the output: it obtains the cost and proceeds backwards calculating the gradients of the hidden units.

Training challenges

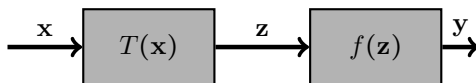
Training neural networks present several challenges:

- **Initialisation:** If the initial weights are zero, back-propagation fails. Initial weight values should be random.
- **Overfitting:** Neural networks can have millions of parameters. Use regularisation and validation-based early stop to avoid overfitting.
- **Non-convexity:** The cost function has multiple local minima. Retrain from different random starting values.
- **Scaling of the inputs:** Range of input values can affect the values of the weights. Standardise to ensure inputs are treated equally.
- **Architecture:** Different architectures suit different problems.

Transfer learning

A neural network that has been successfully trained for a problem A can be reused for a related problem B , for instance:

- We can leave the early stages **unchanged**, becoming a fixed transformation stage $T(\mathbf{x})$.
- We can **retrain** the late stages using new data $f(\mathbf{z})$.



We are in essence **transferring an already learnt transformation** and reusing it for a different problem:

- No need to train $T(\mathbf{x})$ (same parameters for problems A and B).
- The optimal parameters of $f(\mathbf{z})$ for problem B will be close to the ones found for problem A (shorter training time!).

Agenda

Neural networks

What can perceptrons do?

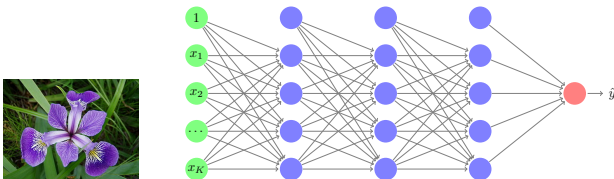
Neural networks as machine learning models

Types of layers

Summary

Fully-connected layer

So far we have considered **fully-connected** (FC) layers are, where all the perceptrons receive all the outputs from the previous layer. FC layers have a **large number of parameters** and training them can be challenging.



Assuming the input of this neural network is an RGB picture consisting of 1000×1000 pixels, how many parameters per perceptron are there in the first FC layer?

Equivariance in grid data

Images and time series are complex data types consisting of individual attributes associated to a **regular grid** defining a **spatial relationship**.

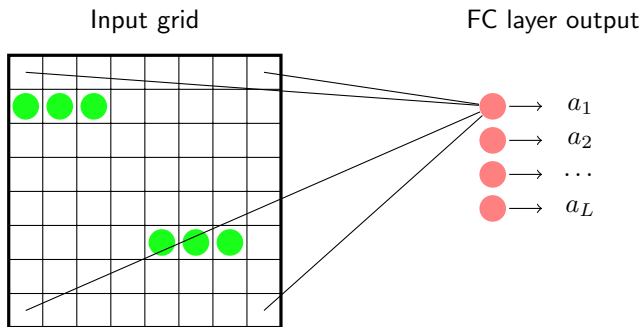
Some grid data exhibit the **equivariance** property, according to which the same pattern can be expected in different locations of the grid.



Equivariance in grid data: FC layer

How would a FC layer identify a short horizontal segment in the following 8×8 grid?

- Each perceptron connected to all inputs: $8 \times 8 + 1$ weights.
- As many perceptrons as potential locations, L .
- Total of $L \times (8 \times 8 + 1)$ weights.

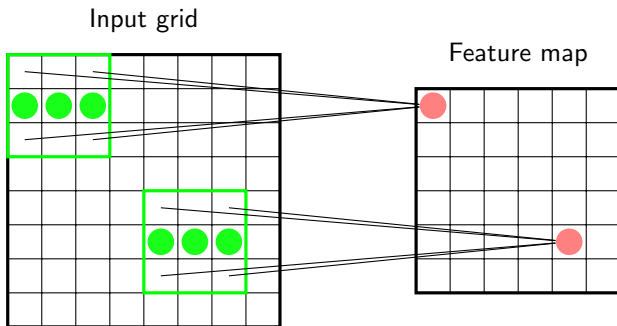


Equivariance in grid data: The convolutional layer

Convolutional layers impose additional restrictions:

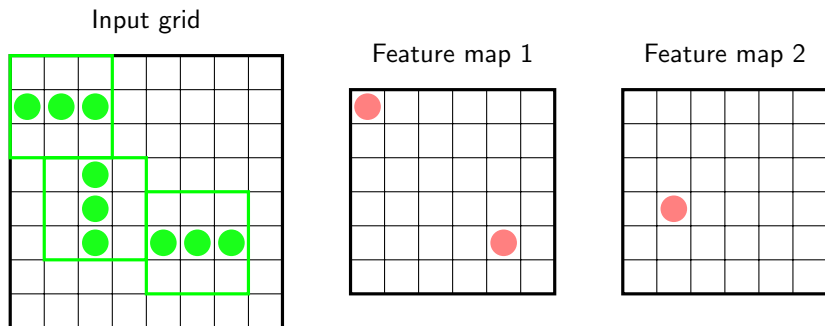
- Perceptrons are arranged as a grid known as **feature map**,
- focus on different **limited regions** in the input grid and
- **share** their parameters, represented as a grid called **kernel**.

The feature map is efficiently calculated as a **convolution** of the kernel and the input or in other words, **filtering** the input with the kernel.



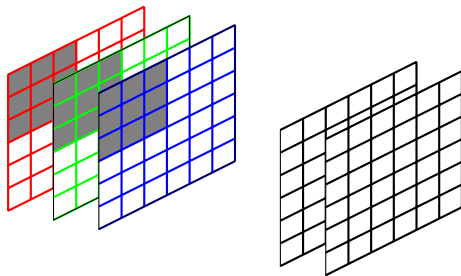
Convolutional layers

Convolutional layers can have **several feature maps**, each of which is associated to a **different concept**. They form a **stack** of maps.



Convolutional layers

- The **dimensions of a kernel** are $H \times W \times D$, there H is the height, W is the width and D is the depth (number of input feature maps).
- The total number of **weights per kernel** is $H \times W \times D + 1$ (bias).
- **Training** a convolutional layer means using data to tune the weights of each kernel.



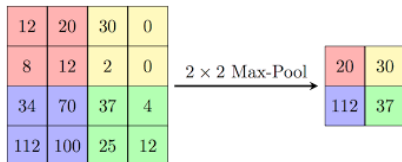
Pooling layers

Pooling **reduces the size** of feature maps. Pooling layers are defined by size of the area they are reducing to a single number and are inserted between successive convolutional layers.

They come in two flavours:

- **Max pooling:** The output is the largest value within the filter area.
- **Average pooling:** The output is the average of the values within the filter area.

Note that pooling layers do not need to be trained!



Deep learning architectures

Deep neural networks are **not arbitrary** sequences of **arbitrary** layers. On the contrary, they have a **predefined architecture** that is suitable for a specific goal.

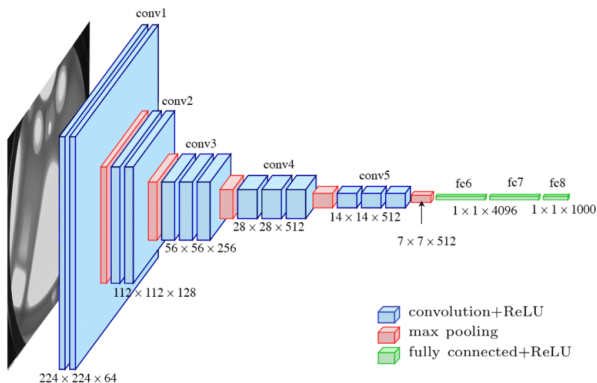
In classification, it is common to see architectures in which:

- The first layers define a **few, simple** concepts, the last layers define **many, complex** concepts.
- Feature maps **shrink** as we move deeper into the network.

The same intermediate concepts can be useful for different goals. We can use **transfer learning** to reuse existing solutions.

Example: The VGG16 network

- VGG16 was designed for the ImageNet Large-Scale Visual Recognition Challenge (dataset of images belonging to 1000 classes).
- VGG16 has 16 layers, 3×3 kernels and 138 million weights.



Agenda

Neural networks

What can perceptrons do?

Neural networks as machine learning models

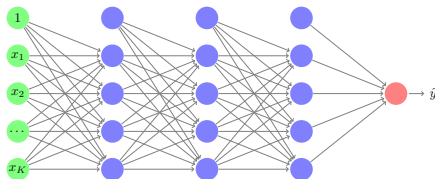
Types of layers

Summary

Neural networks: the three views

We can approach neural networks from three angles:

- **Functional:** system that **maps** an input x to an output \hat{y} . By tuning its set of parameters \mathbf{W} , we change the mapping.
- **Cognitive:** system that **creates new concepts** from raw data and other concepts. By tuning \mathbf{W} , we create different concepts.
- **Computational:** pipeline of interconnected computing units called **perceptrons**. By tuning \mathbf{W} , we change the computation.



Neural networks: A family of machine learning models

Neural networks should not be seen as a machine learning model in the same sense as, for instance, a logistic regression classifier.

- Neural networks are a **family** of machine learning models.
- Each neural network has an **architecture**, the simplest one of which is one single perceptron.
- Some architectures are **not substantially different** from other machine learning models (e.g. the perceptron is a linear classifier).
- We can see neural networks as a **framework** to create new models.
- We **train** one specific architecture and can use **validation** approaches to choose the right architecture.

We should **avoid** saying *neural networks perform well* for a given problem. Instead, we should say *this neural network architecture performs well*.

Neural networks: Universal machines

For any problem, we can find a neural network architecture that solves it. In this sense, neural networks are said to be universal machines.

- Note that this does **not mean** that a specific neural network architecture can solve any problem.
- The existence of a suitable neural network architecture does **not imply** that we will be able to find it.
- Flexibility is achieved adding **complexity**, which increases the risk of **overfitting**.
- **Neural network experts** design the right architecture for a problem and reuse existing solutions using the principles of transfer learning.
- **Neural network brutes** are unaware of the Monkey Theorem and use all the computational power and data available to train as many complex architectures as possible. Their **carbon footprint** is huge.