

QUEEN MARY SCHOOL HAINAN  
QUEEN MARY UNIVERSITY OF LONDON

# QHM5703 Principles of Machine Learning (Introduction to) Neural Networks and Deep Learning

Dr Nikesh Bajaj  
Dr Jesús Requena Carrión

Week 15: 17/18 Dec 2025

# Agenda

Patterns and structure

Grid Data

Neural networks

What can perceptrons do?

Neural networks as machine learning models

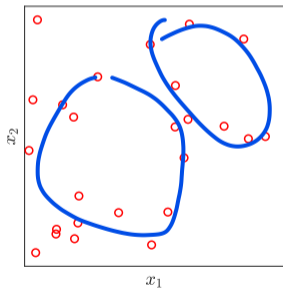
Types of layers

Summary

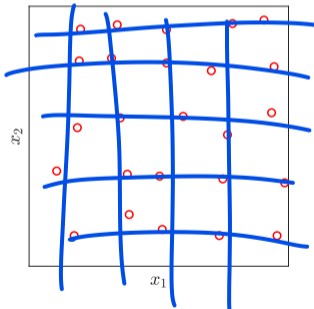
## Patterns and structure

**Patterns** are regularities in our data and **structure** is a regularity in our target population. Machine learning project relies on discovering the underlying structure by identifying patterns in data.

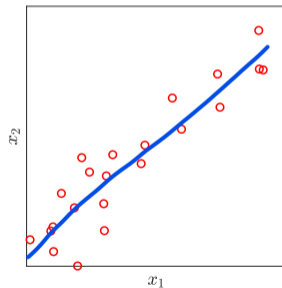
The main challenge is to distinguish **relevant** from **spurious patterns**. Which of the patterns below suggests the **least underlying structure**?



(a)



(b)

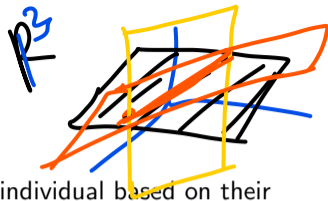
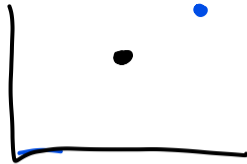


(c)

Go to [menti.com](https://menti.com) and use code: **DMT: 7318 3011** — **ICS: 5841 2762**

## Short and wide datasets

$$N < d$$

 $\mathbb{R}^2$ 

Using the dataset below, train a model that predicts the salary of an individual based on their age, height, weight, gender and postcode.

Short

wide

ID	Age	Height [m]	Weight [kg]	Gender	Postcode	Salary
0	30	1.68	75	M	EC2Y 8DS	100
1	40	1.78	70	F	E14 0QR	200
2	35	1.7	85	M	WC2H 8LH	150

Can you think of patterns in above dataset, to predict the salary from other attributes?

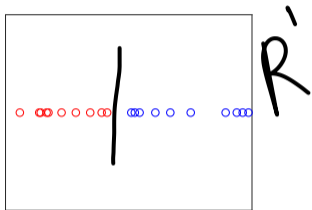
$$N=3$$

 $\mathbb{R}^5$ 

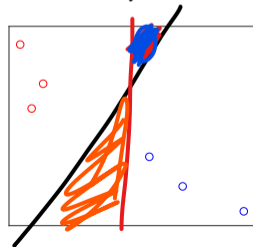
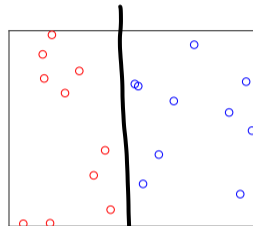
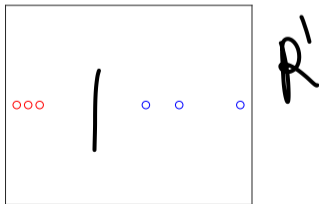
$$N < d$$

# The Curse of Dimensionality

As the dimensionality increases, **data becomes sparser**. What happens if we add irrelevant attributes? What if we have fewer samples?



$n = 6$



# The Curse of Dimensionality

We find high dimensional datasets in many scenarios:

- **Complex** data types, such as **grid data** (pictures, audio files).
- Situations where we have **little prior knowledge** and we record everything, just in case.

Collecting as many attributes as possible seems like a good idea. Is it? Not necessarily: it **gives Randomness more opportunities to fool you.**

The **curse of dimensionality** is a warning. Irrelevant attributes **do not cancel each other out**, they show up as spurious patterns. Adding more attributes (*just in case*) can actually result in worse-performing models.

We can use **feature selection** techniques to reduce the dimensionality of the problem but first, let's use our **domain knowledge**.

Don't let randomness fool you!

# Agenda

Patterns and structure

**Grid Data**

Neural networks

What can perceptrons do?

Neural networks as machine learning models

Types of layers

Summary

## Grid Data

$x_1$	$x_2$	$x_3$	$y$
.	.	.	.
.	.	.	.
.	.	.	.
.	.	.	0

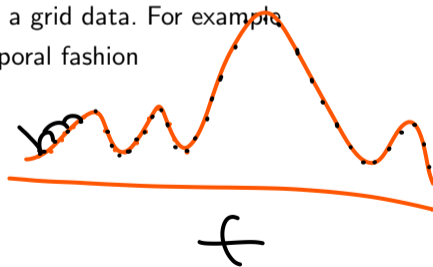


We have seen many datasets, we specify a particular kind of data as **Grid Data:**

**Grid Data:** A data that lives on a grid. The values (predictors) of data are arranged in a way where the notion of neighborhood is used.

For instants, temporal or spatial arrangements of values is a grid data. For example

- Audio/Speech signal - samples are arranged in a temporal fashion
- Image - pixels are arranged in a spatial fashion.
- Graph data - nodes connected to neighbors.



# Agenda

Patterns and structure

Grid Data

**Neural networks**

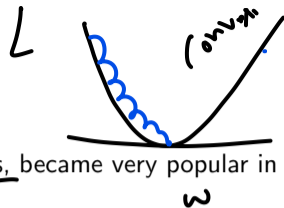
What can perceptrons do?

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.

▪ **Convexity.**

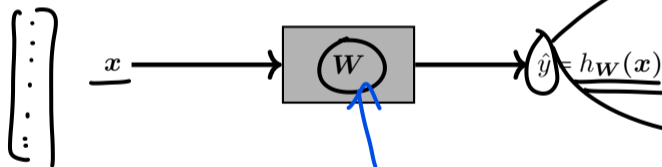
From 2010 onwards there has been a resurgence of interest in neural networks because of:

- 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. 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).



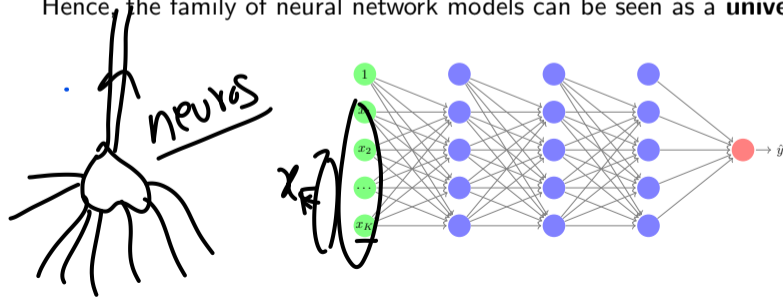
$W$  denotes all the weights of the neural network and  $h_W(x)$  indicates explicitly that the prediction  $\hat{y}$  depends on  $W$ .

# What is a neural network?

Neural networks consist of simple, interconnected computing 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**.

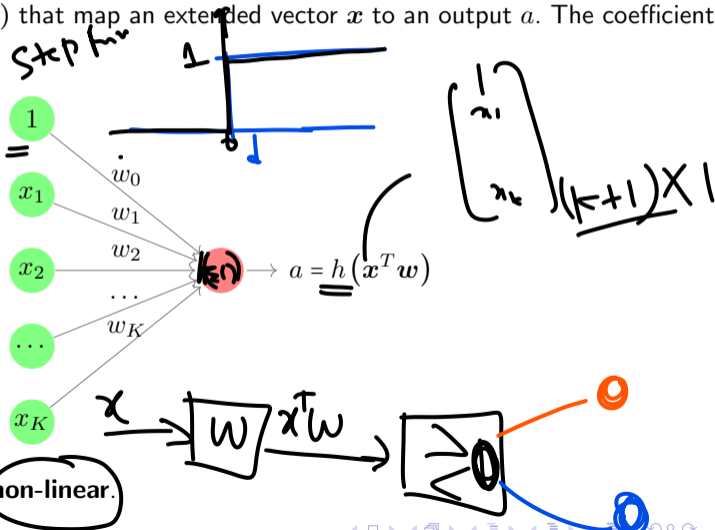


# The perceptron

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**.

logistic key:

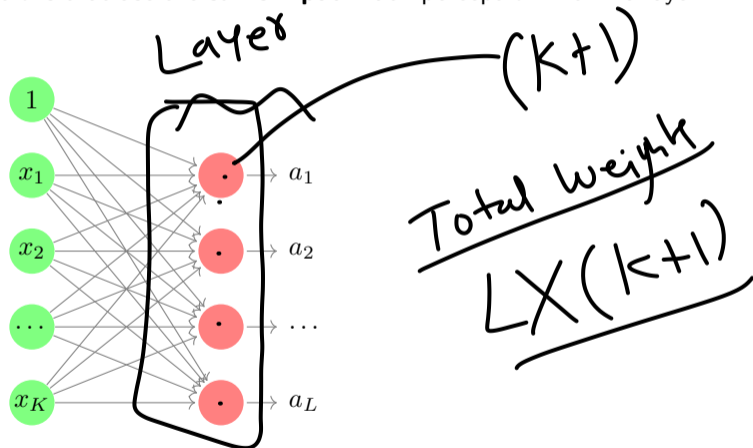
$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_k \end{pmatrix}$$



The activation function is non-linear.

## Layers

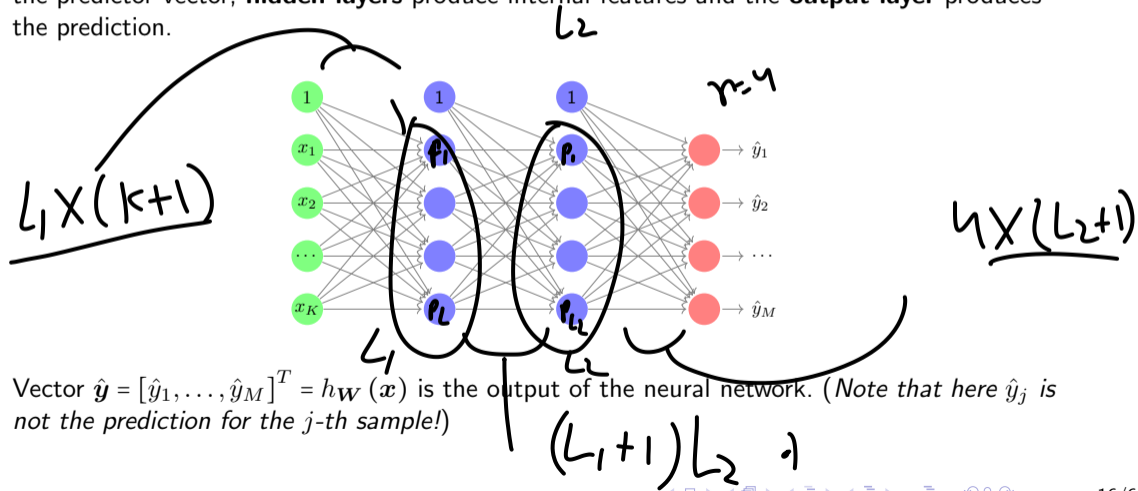
A **layer** is a collection of perceptrons that use the **same input**. Each perceptron within a layer produces a separate output.



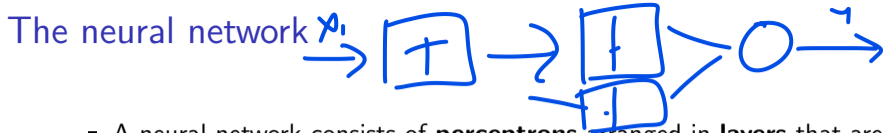
A layer consisting of  $L$  perceptrons and  $K$  inputs has  $L \times (K + 1)$  weights.

# 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.

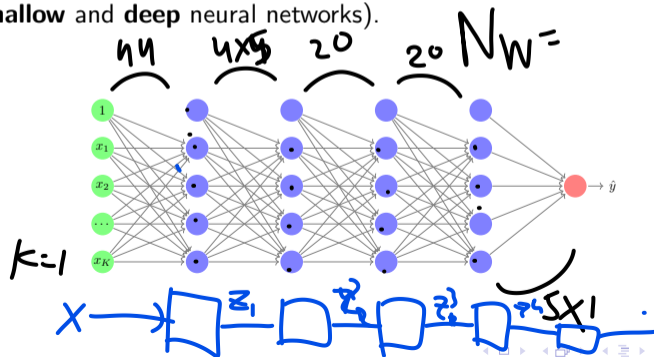


Vector  $\hat{\mathbf{y}} = [\hat{y}_1, \dots, \hat{y}_M]^T = h_{\mathbf{W}}(\mathbf{x})$  is the output of the neural network. (Note that here  $\hat{y}_j$  is not the prediction for the  $j$ -th sample!)



- 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).

pipeline



# Agenda

Patterns and structure

Grid Data

Neural networks

**What can perceptrons do?**

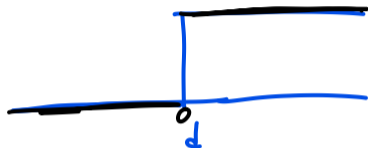
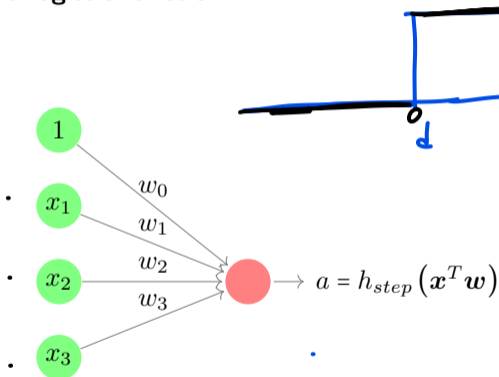
Neural networks as machine learning models

Types of layers

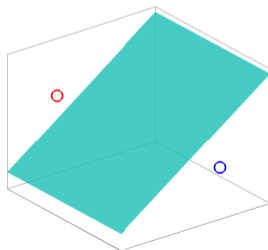
Summary

# 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**?



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



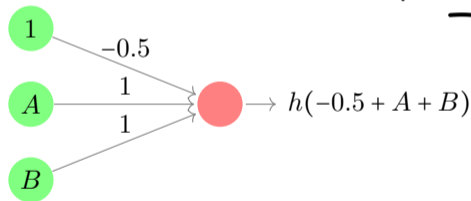
In a 3D predictor space, a perceptron with a step activation function separates two decision regions by a plane!



# The perceptron as a basic logical function

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

$h \rightarrow$  Step  $f_n$



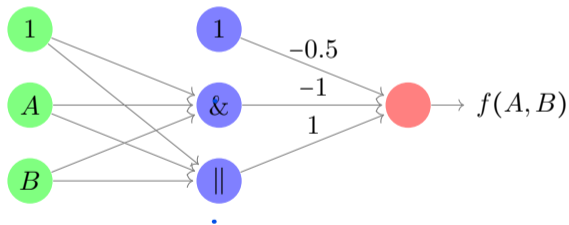
$A$	$B$	$f(A, B)$
0	0	·
0	1	·
1	0	·
1	1	·

$|| \leftarrow$  OR gate

## Adding depth to derive new logical functions

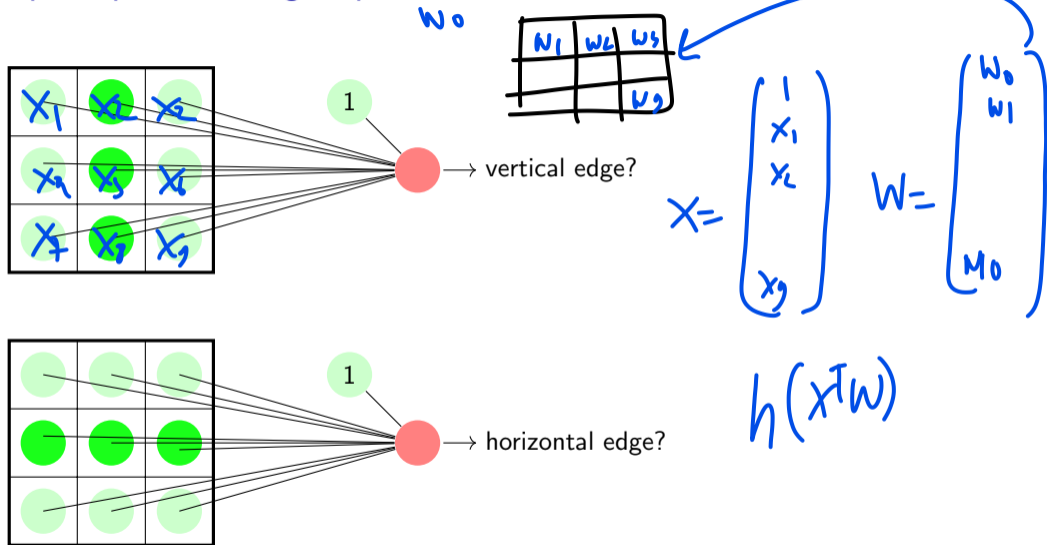
XOR

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



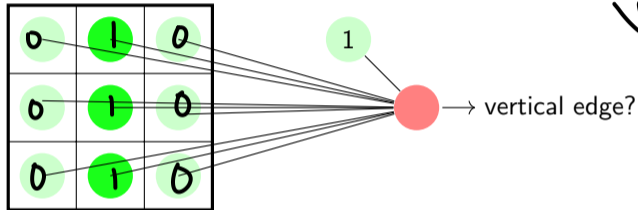
$A$	$B$	&		$f$
0	0			
0	1			
1	0			
1	1			

# The perceptron as a grid pattern detector

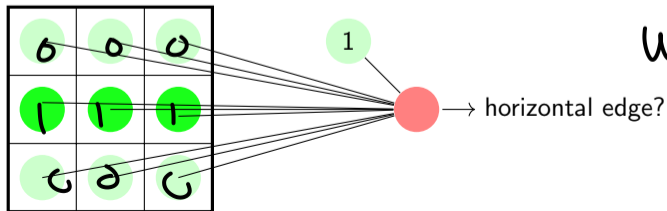


# The perceptron as a grid pattern detector: Question

If dark green circles have high value ( $> 0$ ) and light green circles have low values ( $< 0$ ), what should be the weights of following perceptrons?

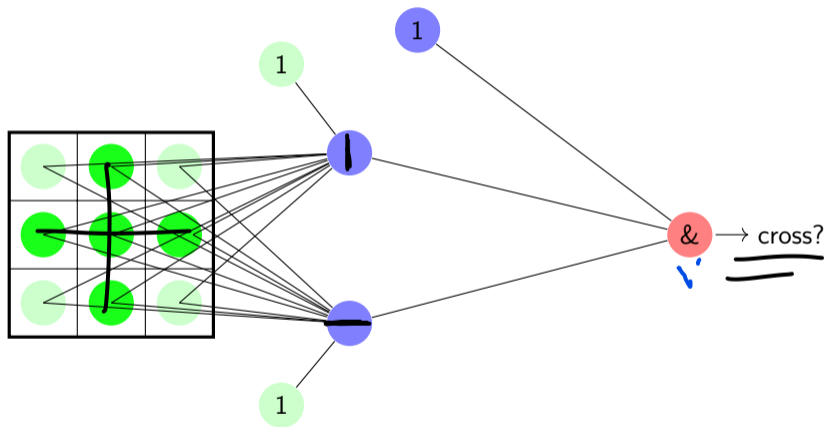


$$W = \begin{bmatrix} -1 & 1 & -1 \\ -1 & 1 & -1 \\ -1 & 1 & -1 \end{bmatrix}$$



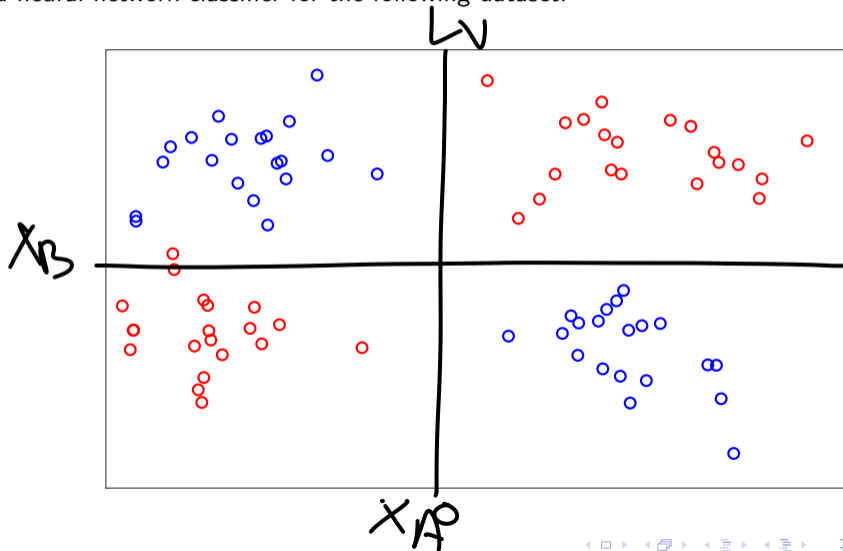
$$W = \begin{bmatrix} -1 & -1 & -1 \\ 1 & 1 & 1 \\ -1 & -1 & -1 \end{bmatrix}$$

# Adding depth to detect derived grid patterns



# Combining linear classifiers and logical functions: Question

Build a neural network classifier for the following dataset.



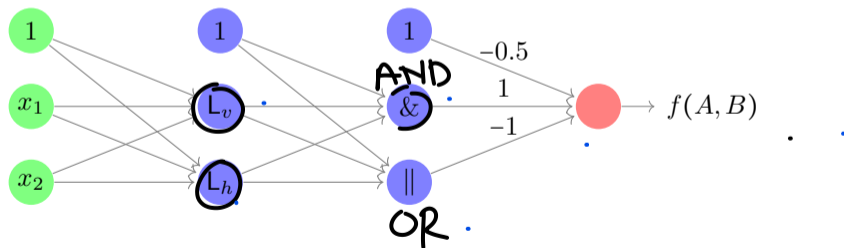
$$x = \begin{bmatrix} 1 \\ x_A \\ x_B \end{bmatrix}$$

$$L_v = [0 \ 1 \ 0]$$

$$L_h = [0 \ 0]$$

## 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!

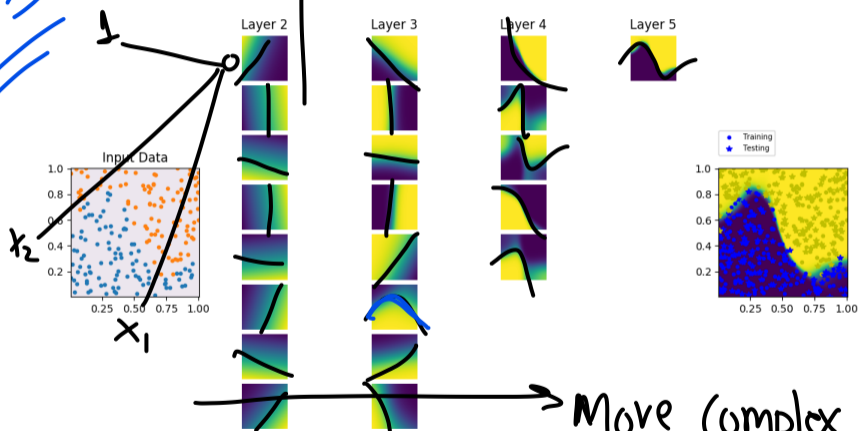
# Complex Features: Example

#W = ?

New concept

How many weights?

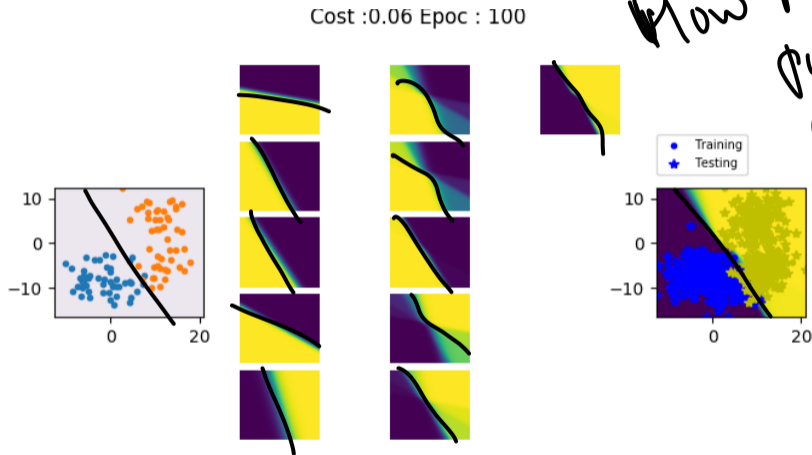
Cost : 0.07 Epoc : 600



Move complex features

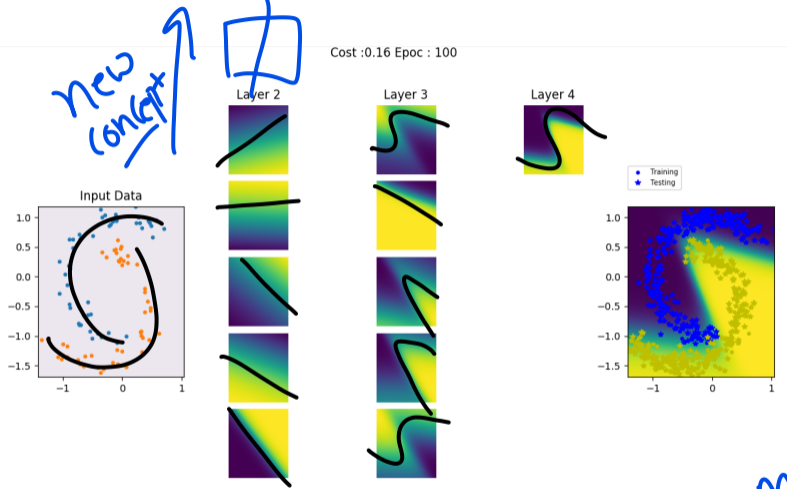
Source: [https://github.com/Nikeshbajaj/DeepLearning\\_from\\_scratch](https://github.com/Nikeshbajaj/DeepLearning_from_scratch)

## Complex Features: Example 2



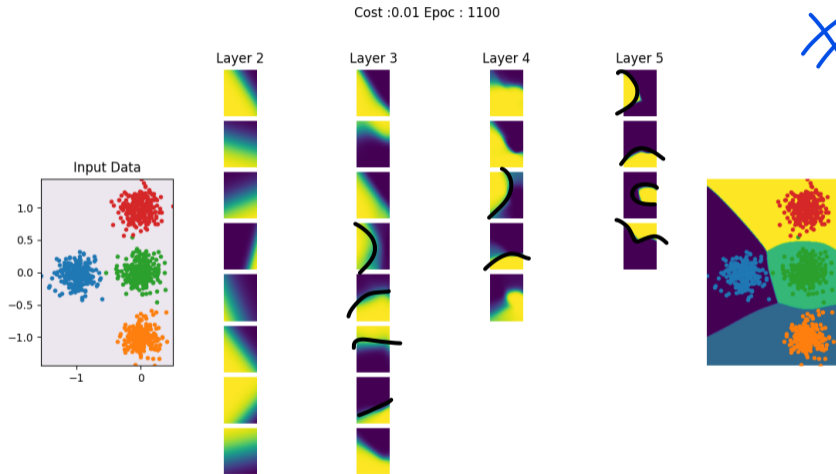
Source: [https://github.com/Nikeshbajaj/DeepLearning\\_from\\_scratch](https://github.com/Nikeshbajaj/DeepLearning_from_scratch)

# Complex Features: Example 3



Source: [https://github.com/Nikeshbajaj/DeepLearning\\_from\\_scratch](https://github.com/Nikeshbajaj/DeepLearning_from_scratch)

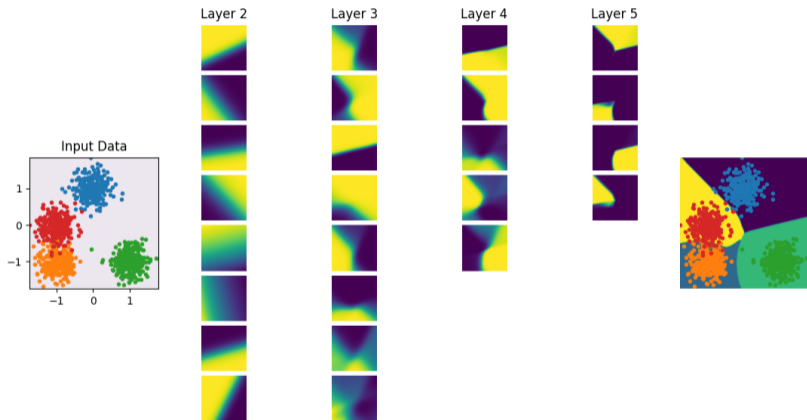
# Complex Features: Example 4



Source: [https://github.com/Nikeshbajaj/DeepLearning\\_from\\_scratch](https://github.com/Nikeshbajaj/DeepLearning_from_scratch)

# Complex Features: Example 5

Cost :0.04 Epoc : 300



Source: [https://github.com/Nikeshbajaj/DeepLearning\\_from\\_scratch](https://github.com/Nikeshbajaj/DeepLearning_from_scratch)

Visualisation: <https://nikeshbajaj.github.io/teaching/ml101/cover.gif>

# Agenda

Patterns and structure

Grid Data

Neural networks

What can perceptrons do?

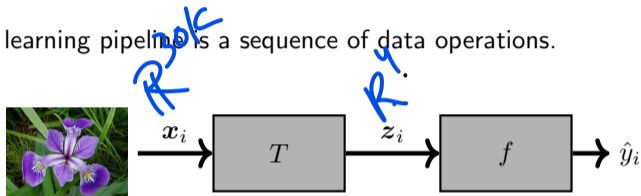
**Neural networks as machine learning models**

Types of layers

Summary

## Machine learning pipelines: Reminder

A machine learning pipeline is a sequence of data operations.



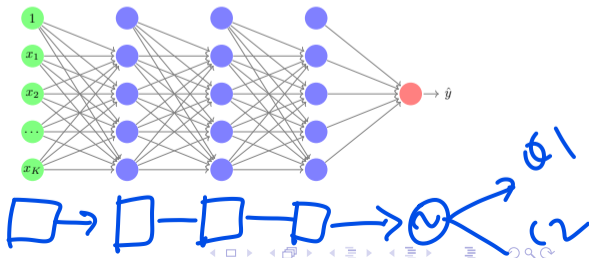
A simple pipeline consists of a first transformation stage followed by a machine learning model:

- 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:

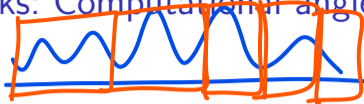
- Each perceptron defines one **derived feature** or **concept** using other features, raw or derived.
- **Increasing the number of perceptrons** in a layer allows to create more new concepts per layer.
- **Increasing the number of layers** (deeper networks) allows to create concepts of increasing complexity.



# Deep neural networks: Computational angle

*Train*

①

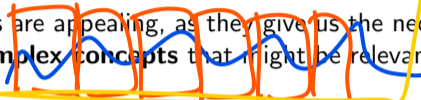


305

18

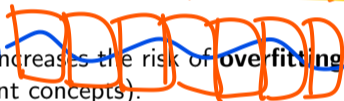
Large neural networks are appealing, as they give us the necessary **flexibility to create new and increasingly complex concepts** that might be relevant to make a prediction.

②



However:

③



- 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.

*Train* → 10  
*Val* → 8

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  $l(\mathbf{W})$ .

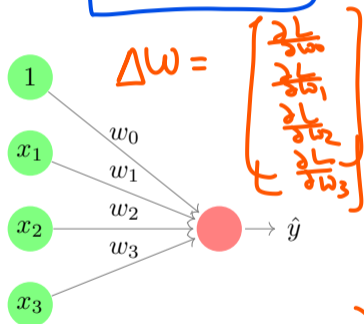
Obtaining the gradient is easy, but can be computationally expensive. **Back-propagation** is an efficient algorithm to **compute the gradient**. This gradient is then used by the optimisation algorithm to update  $\mathbf{W}$ .

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.

# Gradient descent and back-propagation

$$L = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

Linear:  $\hat{y} = w_0 + \sum_i w_i x_i$



$$\frac{\partial L}{\partial w_1} = \frac{1}{N} \sum_{i=1}^N 2(\hat{y}_i - y_i) \cdot \frac{\partial \hat{y}}{\partial w_1}$$
$$\frac{\partial L}{\partial w_1} = \frac{1}{N} \sum_{i=1}^N (2(\hat{y}_i - y_i) \cdot x_1)$$

$$w_1 = w_1 - \delta \frac{\partial L}{\partial w_1}$$

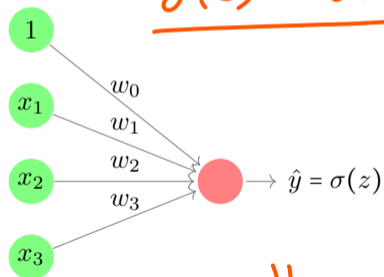
$$w = w - \delta \Delta W$$

$$\frac{\partial L}{\partial w_0}, \frac{\partial L}{\partial w_2}, \frac{\partial L}{\partial w_3}$$

# Gradient descent and back-propagation

Non-Linear:  $\hat{y} = \sigma(w_0 + \sum_i w_i x_i)$   
 $\sigma(z) = \frac{1}{1+e^{-z}}$

$$\sigma'(z) = \sigma(z)(1-\sigma(z))$$



$$\begin{aligned} \frac{\partial L}{\partial w_1} &= \frac{\partial}{\partial w_1} (y - \hat{y}_i)^2 \\ &= -2(y - \hat{y}_i) \frac{\partial \hat{y}}{\partial w_1} \\ &= -2(y - \hat{y}) \sigma(z) (1 - \sigma(z)) \\ &\quad \cdot \frac{\partial}{\partial w_1} (w_0 + \sum w_i x_i) \end{aligned}$$

$$\frac{\partial L}{\partial w_1} = -2(y - \hat{y}) \sigma(z) (1 - \sigma(z)) x_1$$

# Gradient descent and back-propagation

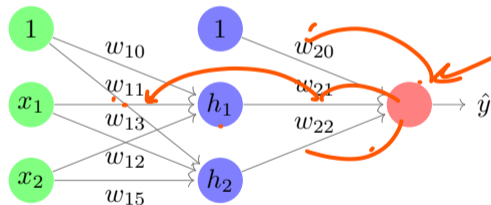
$$\hat{y} = \sigma(w_{20} + \sum_i w_{2i} h_i)$$

$$h_1 = \sigma(w_{1,0} + \sum_i w_{1i} x_i)$$

$$h_2 = \sigma(w_{13} + \sum_i w_{1i+3} x_i)$$

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

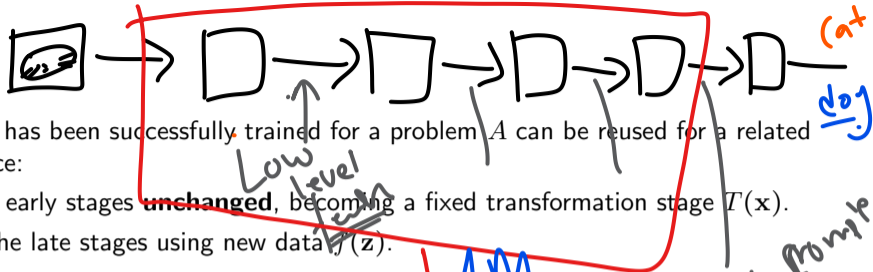
$$\begin{aligned} \frac{\partial L}{\partial w_{11}} &= \frac{\partial (y - \hat{y})^2}{\partial w_{11}} \\ &= 2(y - \hat{y}) \cdot \frac{\partial \hat{y}}{\partial w_{11}} \end{aligned}$$



## Training considerations

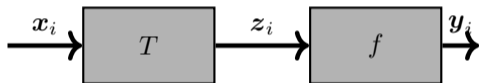
- **Initialisation**: If the initial weights are zero, back-propagation fails. Initial weight values should be random.
- **Overfitting**: Neural networks can have millions of parameters. Use regularisation and validation-based early stop to avoid overfitting.
- **Non-convexity**: The cost function has 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  $T(x)$ .
- We can **retrain** the late stages using new data  $f(z)$ .



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

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



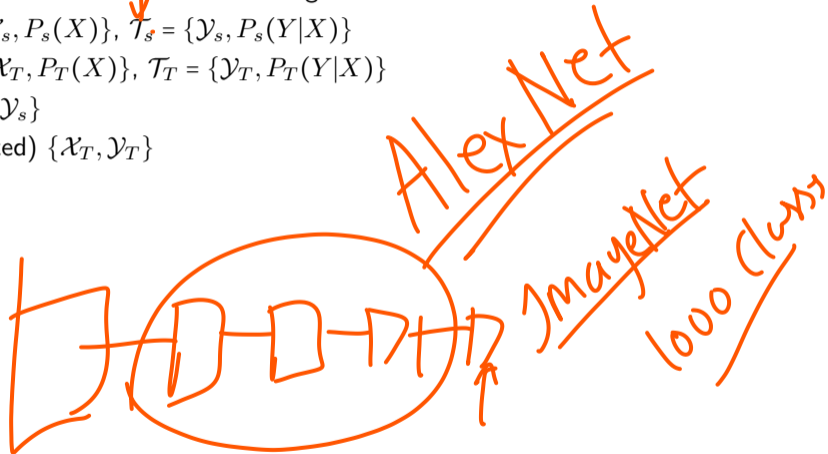
# Transfer learning: A little more with math

*cat/dog*

Key mathematical concepts to define transfer learning

- Source:  $\mathcal{X}_s, \mathcal{D}_s = \{\mathcal{X}_s, P_s(X)\}, \mathcal{T}_s = \{\mathcal{Y}_s, P_s(Y|X)\}$
- Target:  $\mathcal{X}_T, \mathcal{D}_T = \{\mathcal{X}_T, P_T(X)\}, \mathcal{T}_T = \{\mathcal{Y}_T, P_T(Y|X)\}$
- Source dataset  $\{\mathcal{X}_s, \mathcal{Y}_s\}$
- Target dataset (limited)  $\{\mathcal{X}_T, \mathcal{Y}_T\}$

- $\mathcal{D}_s \neq \mathcal{D}_T$
- $\mathcal{T}_s \neq \mathcal{T}_T$



# Transfer learning: Examples

Examples: similarity of domain

- Image datasets ✓
- Audio datasets ✓
- Biomedical signals ✓

# Agenda

Patterns and structure

Grid Data

Neural networks

What can perceptrons do?

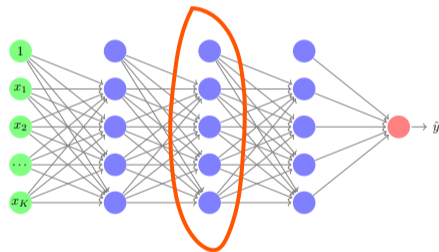
Neural networks as machine learning models

**Types of layers**

Summary

## Fully-connected layer

So far we have considered **fully-connected** (FC) layers, 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  $1000 \times 1000$  pixels, how many parameters per perceptron are there in the first FC layer?

## Equivariance in grid data

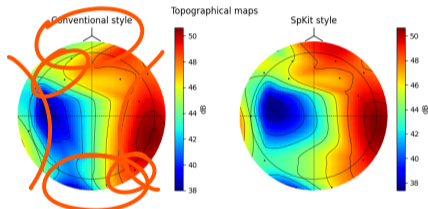
Images and time series are complex data types consisting of individual attributes associated to a **regular grid** defining a spatial relationship.

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

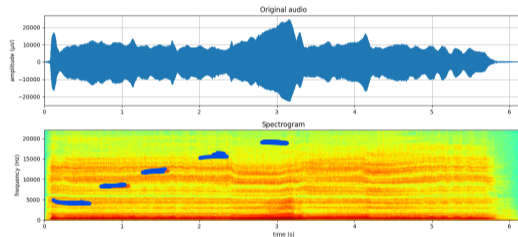


# Equivariance in grid data: Examples

Images with no Equivariance.



Brain Activity



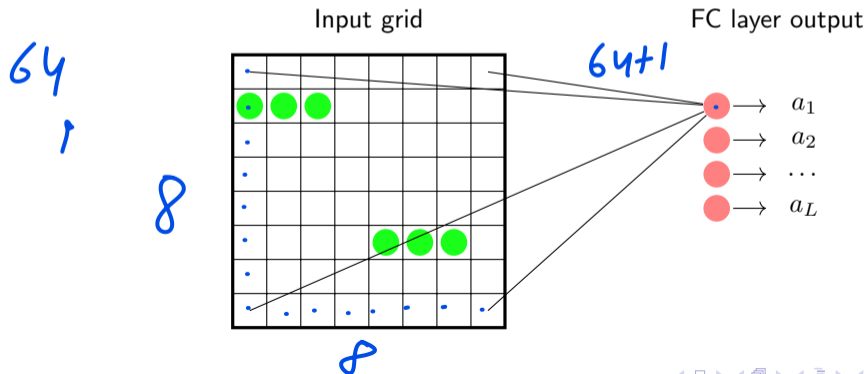
Spectrogram of an audio

MRI Images  
CT Scan Images

## Equivariance in grid data: FC layer

How would a FC layer identify a short horizontal segment in the following  $8 \times 8$  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.

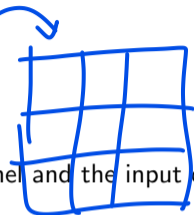


# Equivariance in grid data: The convolutional layer

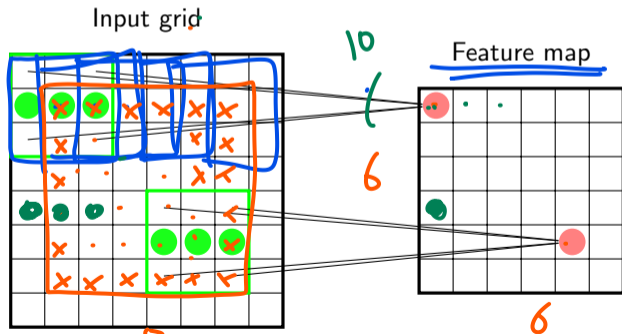
Convolutional layers impose additional restrictions:

- 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.



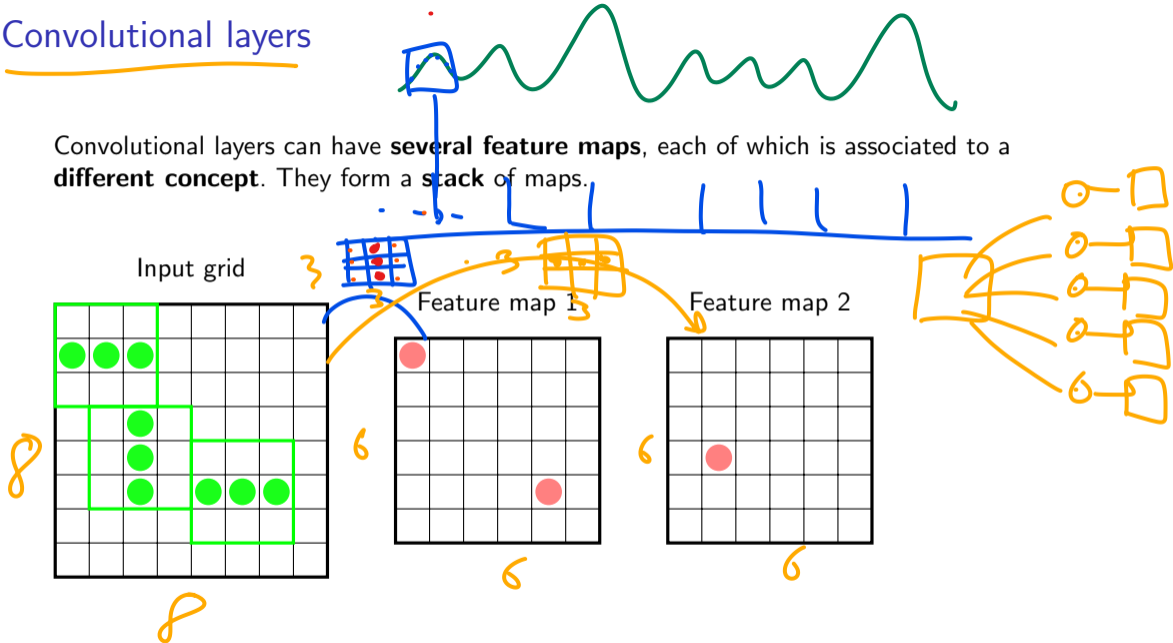
Convolution



$$\begin{array}{r} 36 \times 10 \\ 360 \\ \hline 10 \end{array}$$

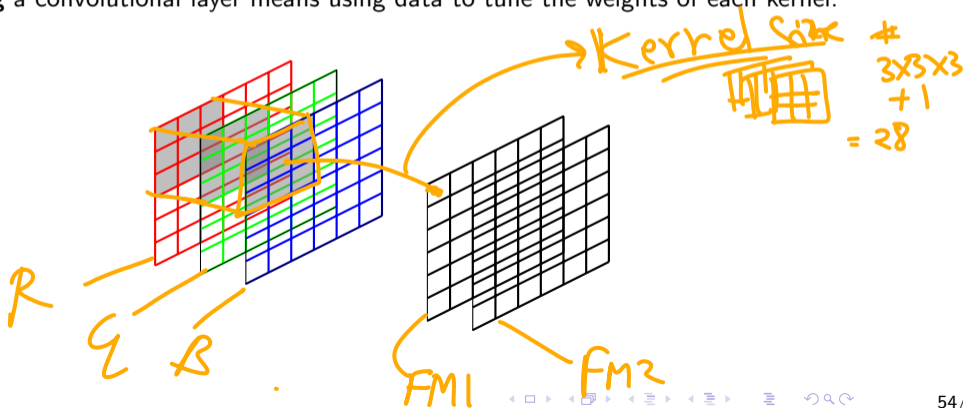
# Convolutional layers

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  $H \times W \times D$ , where  $H$  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.



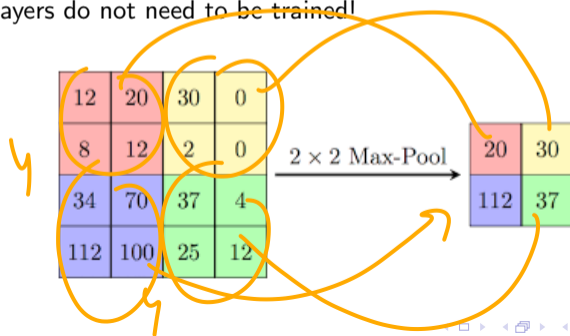
## Pooling layers

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.

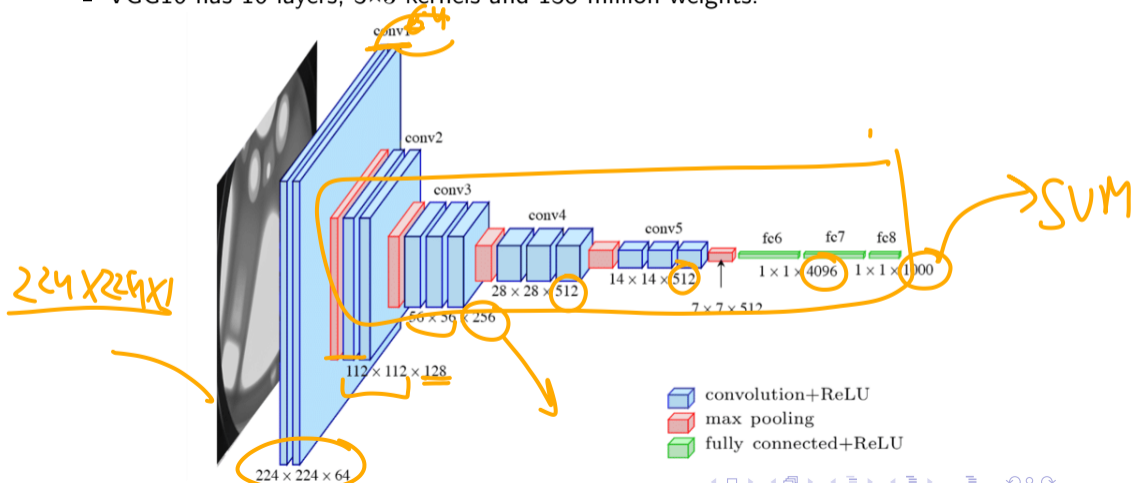
In classification, it is common to see architectures in which:

- 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.

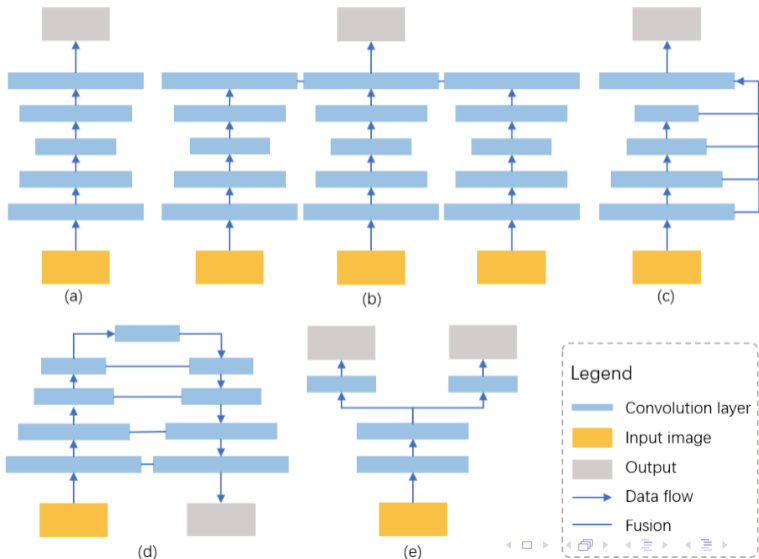
The same intermediate concepts can be useful for different goals. We can use **transfer learning** to reuse existing solutions.

# 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,  $3 \times 3$  kernels and 138 million weights.



# More architectures



# Agenda

Patterns and structure

Grid Data

Neural networks

What can perceptrons do?

Neural networks as machine learning models

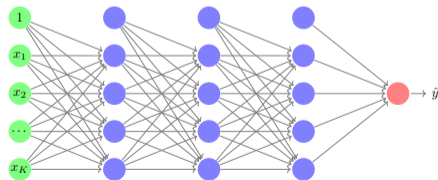
Types of layers

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  $\mathcal{W}$ , we change the mapping.
- **Cognitive**: system that **creates new concepts** from raw data and other concepts. By tuning  $\mathcal{W}$ , we create different concepts.
- **Computational**: pipeline of interconnected computing units called **perceptrons**. By tuning  $\mathcal{W}$ , 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.

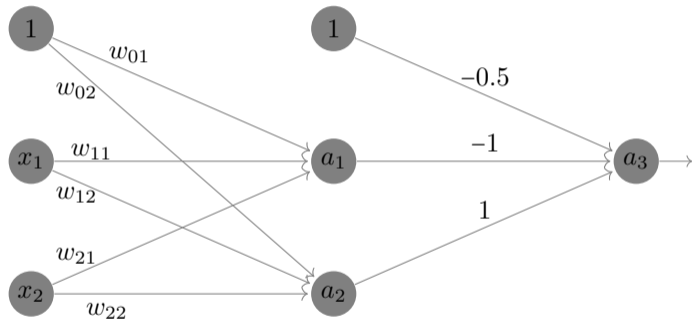
We should **avoid** saying *neural networks perform well* for a given problem. Instead, we should say *this neural network architecture performs well*.

## 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.

# Ex 1





Queen Mary  
University of London