

(Human) Network Analysis

Lecture 1: **Principles of deep learning in artificial networks**

Ben Harvey

1

Welcome to Network Analysis. This is a course from the faculty of social sciences, where we will focus on two types of human-based networks.

In the first half of the course we will look at deep convolutional neural networks, a type of machine learning network whose design is inspired by the function of neurons in the brain, and the structure of networks among these neurons. We will look at how these artificial and biological networks are related.

In the second half of the course, my colleague Jiamin Ou from the department of Sociology will teach you about social networks, both online and in human society. She will focus on how information transfers between human in these networks.

In deep networks, the network units are neurons, and very large sets of neurons interacting together form networks with emergent properties: the network does things that any single neuron can't do.

The largest scale network here is a complete human brain, whose emergent properties underlie human behaviour.

In social networks the network unit is a behaving human, so the network in the deep networks part of the course is the unit in the social network part of the course.

Ideally, we would model this human-sized network unit as a network in itself, and model every neuron's activity to predict human behaviour. There have been great advances in deep learning in recent years, but we are not yet at the stage where we can make a complete model of a behaving human, let alone use deep learning models to simulate several interacting humans.

Neural network models should be usable to model the agents in social network models, but humans are just too complicated for our current computing power.

As a result, there is a disconnection between the first half and the second half of the course. In social networks, we have to use very simplified models of

human behaviour to allow the models to scale up to simulate large social networks.

First, let's go through an outline of the course so you know what to expect and what we expect from you.

Course goals

- Explore the relationship between cognitive science and AI
- Focus on deep learning in artificial machine learning networks and comparison to biological systems
 - Which biological processes do deep networks imitate?
 - What is missing in artificial networks?
 - What might make AI/machine learning more like biological intelligence/learning
- Become familiar with the use of AI in cognitive science research
- Build some deep learning networks to do human-like tasks

2

Some of this lecture content is complex, and you have very different backgrounds.

But it will be examined, and you must pass this exam to pass the course.

So we will go slowly, start from first principles and avoid assuming any knowledge.

Please ask questions when you don't understand: the explanations and discussion that result are very valuable to the class format.

For remote teaching, we find questions work best as text messages in the Teams chat

Assessment

- In groups of 4, students will complete two lab assignments to build simple deep learning systems to solve computer vision and linguistics tasks. Students will be graded on two written reports of their work. Grades depend on depth and completion (20% each).
- Each student will complete an individual assignment related to each lab assignment (10% each)
- Students' understanding of lectures and reading assignments will be assessed in a final exam that determines 40% of the final grade.
- You are required to average a passing grade (5.5) **across the exam and individual assignments** to pass the course. Students scoring between 4.0 and 5.5 qualify for a repair exam.

3

Date	Time	Format	Teacher	Room
10/02	11:00-12:45	Lecture 1	Harvey	Remote (Teams)
10/02	15:15-17:00	Lab 1	Overvliet	Remote (Teams)
12/02	11:00-12:45	Lab 1	Overvliet	Remote (Teams)
17/02	11:00-12:45	Lecture 2	Harvey	Remote (Teams)
17/02	15:15-17:00	Lab 1	Overvliet	Remote (Teams)
19/02	11:00-12:45	Lab 1	Overvliet	Remote (Teams)
24/02	11:00-12:45	Lecture 3	Harvey	Remote (Teams)
24/02	15:15-17:00	Lab 1	Overvliet	Remote (Teams)
26/02	11:00-12:45	Lab 1	Overvliet	Remote (Teams)
03/03	11:00-12:45	Lecture 4	Harvey	Remote (Teams)
03/03	15:15-17:00	Lab 1	Overvliet	Remote (Teams)
05/03	11:00-12:45	Lab 1	Overvliet	Remote (Teams)
07/03	23:59	Deadline Lab 1: Group Assignment		
09/03	23:59	Deadline Individual Assignment 1		
16/04	15:15-17:15	EXAM	Remindo (remote)	
?/06	?	Repair exam	Remindo (remote)	

4

Here you can see the schedule for the first few weeks.

You have classes on Fridays where you can work together and get help with your lab assignments, while Wednesday classes will be lectures from me.

This leaves you with 16 classroom hours when you can get help with the assignment. This should be enough if you are working efficiently, but may not be.

If not, we expect you to work on this between classes too, which requires arranging a time with your group when you can work together. Outside class time, you won't be able to get help from your lab teachers, or show them your work.

So if you are falling behind, make sure your work is almost complete by the last lab class so you can talk to your teacher in that class. We'll introduce this assignment at the end of this lecture.

Why deep learning?

- AI has made great advances in tasks that are:
 - Described by formal mathematical rules
 - Relatively simple for computers
 - Difficult for humans
- AI had been less effective in tasks that are:
 - Hard/impossible to describe using formal mathematical rules
 - BUT easy for humans to perform
 - Intuitive or automatic
 - Simulation of neural computation

5

Deep learning tasks



Image processing
Game opponents in complex games
Natural language processing
Simulation of biological neural systems

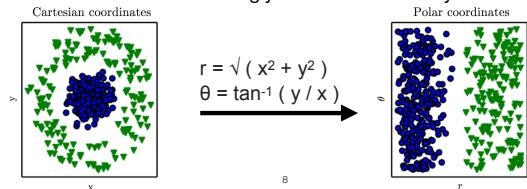
Deep learning approach

- Learn from experience (machine learning)
 - No formal rules of transformations
 - No ‘knowledge base’
 - No logical inference
- Process inputs through a hierarchy of concepts
 - Each concept defined by its relationship to simpler concepts
 - So, build complicated concepts out of simpler concepts

This is exactly what a human is doing when learning about the world
SO the main inspiration used is the brain

Representations & features

- Machine learning performance depends on the **representation** of the case to be classified
 - What information the computer is given about the situation
- Each piece of input information is known as a **feature**
 - The same feature can be represented in different formats
 - Often easy to convert between formats
 - The chosen format strongly affects the difficulty of the task



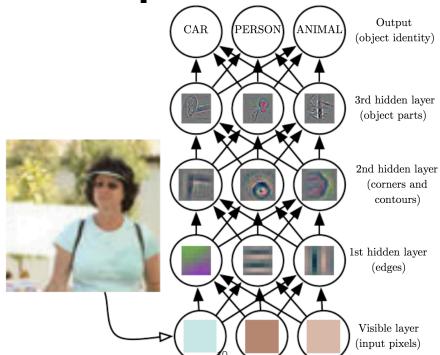
A simple task like this can be solved by choosing the right set of features, with minimal learning necessary.
However, for many tasks, it is hard to know which features or formats of the input are important in determining the output.
And these may be high-level features that need to be extracted first.
Deep learning aims to extract important features, and determine which features are important, through experience

Representations in deep networks

- Useful features may need to be transformed or extracted first
- So deep networks have multiple representations
 - Each is built from an earlier representation
- This can:
 - Transform features to a different format before learning their links to the output
 - Extract complex features from simpler features
- Essentially multiple steps in a program
 - Each layer can be seen as the computer's memory state after executing a set of instructions
 - Deeper networks execute more instructions in sequence
- Just like a computer program, the individual steps are generally very simple
 - Complex outcomes emerge from interactions between many simple steps

9

Representations in deep networks



Here we can see how this abstract description might work in an oversimplified example of object recognition.

The first layer just takes the colour of each pixel.

This is transformed to the edge representation in the next layer by learning common relationships between these pixels.

The edges are then transformed to corners and contours by learning relationships between the edges.

The next layer finds object parts by learning common patterns of corners and contours.

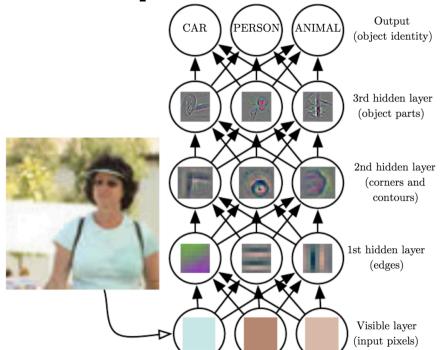
These object parts are then transformed into whole object representations by learning which patterns of object parts correspond to which object type.

We will return to the example of object recognition many times

It's an excellent example of a process that is intuitive and automatic, but hard to formalise or program.

It is also very useful for computers to do, so we can find images on the internet based on their content without a human labelling all this content.

Representations in deep networks



A quick note on notes.

All of these slides will be available online, but you will find I use little text on my slides. This works better in class, but is hard to study from.

I deal with this in two ways. First, my online slides also contain notes with a fairly complete description of what I say. This is very easy to make notes on and study from. Second, I record these lectures and also put these movies on Blackboard.

Please note that your other lecturer won't give you extensive notes, but will also make recordings.

What is a deep network?

- A learning network that **transforms** or **extracts** features using:
 - Multiple **nonlinear** processing units
 - Arranged in multiple **layers** with
 - **Hierarchical organisation**
 - Different levels of **representation** and abstraction

Note that this definition does not specify 'machine' learning.

In this course, we will also look at biological neural networks like the brain, which are also deep networks

Lectures 1-4

- Lecture 1: Principles of deep learning in artificial networks
- Lecture 2: Deep learning in biological neurons and networks
- Lecture 3: Feedforward and recurrent visual processing
- Lecture 4: Simulating biological neural systems

So, we're going to start by looking at artificial networks, a machine learning system. We will compare these to the human brain's biological neural networks these artificial networks aim to imitate.

In my last two lectures, we will look at the stages of visual processing involved in object recognition and other visual tasks in both systems.

This will show what is missing in current artificial networks and how researchers are starting to build deep networks that fill these gaps.

To begin, let's look at what machine learning is, starting artificial networks from the start.

What is machine learning?

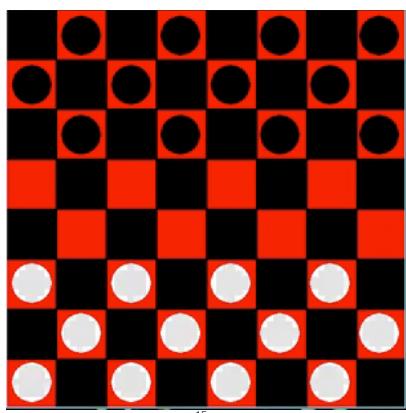
- The field of science that 'gives computers the ability to learn without being explicitly programmed' (Arthur Samuel, 1959)



14

So this field has been around for a long time

Samuel showed a machine learning algorithm many games of checkers (draughts) to learn which board positions are likely to lead to a win. The only programmed rule was the way the pieces are allowed to move.



15

Here we see the black player (a machine learning program) consistently beating the white player (which moves randomly)

What is machine learning?

- The field of science that 'gives computers the ability to learn without being explicitly programmed' (Arthur Samuel, 1959)
- 'A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T , as measured by P , improves with experience E .' (Tom Mitchell, 1998)
- Expressed in organisational terms, not cognitive terms

16

But Samuel's definition of machine learning uses the word 'learn' without any definitions

It also uses 'computer', which is a synonym of 'machine'

Mitchell provided a more formal definition
By describing the fundamental operation in terms of inputs and outputs, this definition avoids suggesting the machine can think.

Alan Turing: "Can machines think?" -
>"Can machines do what we (as thinking entities) can do?"

Learning to filter spam

- 'A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T , as measured by P , improves with experience E .' (Tom Mitchell, 1998)
- **Task:** 'Classify which emails are wanted (not spam) vs unwanted (spam)'
- **Experience:** Watching humans labels emails (training set)
- **Performance:** The proportion of new emails (test set) classified correctly

17

Note that this is a relatively simple machine learning example, and can be done without deep learning.

So, now that we have an idea of what machine learning is, how does machine learning work?

What is a deep network?

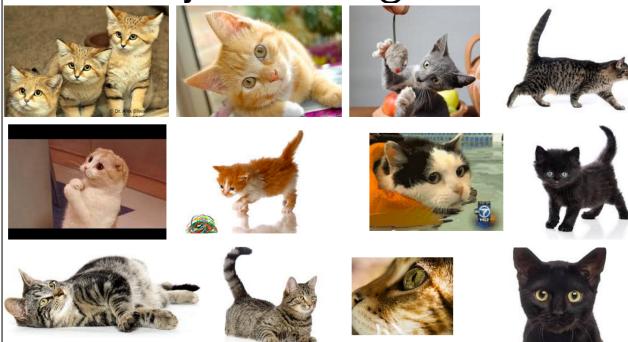
- A learning network that **transforms** or **extracts** features using:
 - Multiple **nonlinear** processing units
 - Arranged in multiple **layers** with
 - **Hierarchical organisation**
 - Different levels of **representation** and abstraction

18

This is a very broad definition, so we will use an example to see what it looks like
The example we will use in the first half of the course is object recognition.

This has been a major goal for deep learning in recent years, and is now largely solved, so we can investigate in depth how this works.
Object recognition may sound like an easy problem for computer vision, but...

Object recognition



Why is it so difficult?

19

The identity of any object has little relationship to its impression on the retina.

Here we see the result of a google image search for pictures of cats. This result is achieved using Google's artificial deep network trained for object recognition.

For this network and also for human vision, objects can be recognised from different viewpoint and sizes, in different positions, and with different lighting conditions

Also, examples of the same class of object often look very different.

For example, you might think that pointed ears tell you this is a cat, but some examples don't have them while clearly being recognised as cats.

So we can't recognise an object directly from its impression on the eye or camera sensor

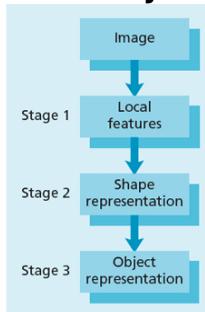
Object selective responses



However, we do know that the brain does this. We are able to recognise the same object across viewpoints, and also the brain contains neurons that respond to specific object types.

A common model for human object-selective responses is face processing

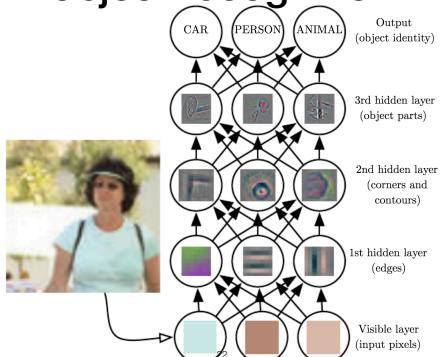
The 20th century view of object recognition



- Stage 1 builds a representation of local image features.
- Stage 2 builds a representation of larger-scale shapes and surfaces.
- Stage 3 matches shapes and surfaces with stored object representations-recognition.

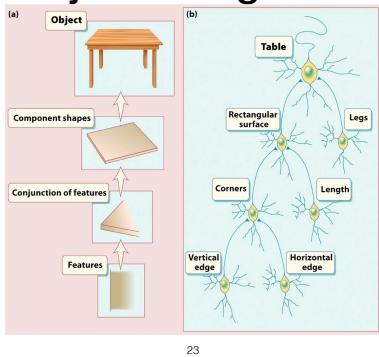
Note this is already a multi-layer, hierarchical approach

The 20th century view of object recognition



This was essentially the example we looked at earlier, here with the input at the bottom.

The 20th century view of object recognition



23

So we might think that our object representations are built from combinations of feature representations. Indeed, many objects are built from parts, so simplified parts that we recognise from all angles might let us build an object viewpoint-independent object representation.

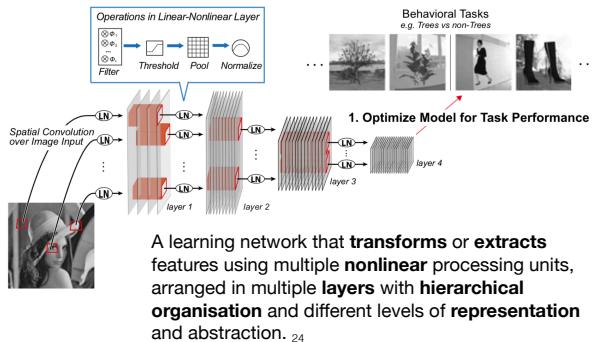
However, no one has ever made a program that can do this for a large set of different objects, and no one has ever found neurons responding to feature conjunctions or component shapes or object parts.

It seems that the features considered in a model like this are too human.

When is a feature a corner and when is it a curve? Is there something in between? Can we define a corner, edge or surface so rigidly?

Essentially, all of these steps tend to limit the network to recognise specific examples, rather than generalise to all possible tables, which is the goal here.

A deep network for object recognition



So we can see this network fulfils all the criteria of a deep network.

It takes an input image and transforms its features to extract the class of object the image contains.

It is arranged in multiple layers, with one feeding into the next, forming a hierarchy. This is all much like the 20th-century idea.

The first layer represents the image pixels, with minimal abstraction, while the last layer captures object identity, which is highly abstract for a computer system. But what happens to get from one to the other is very different from the 20th-century view.

Essentially the difference is that the middle network layers do not respond to concrete concepts like corners and object parts: they respond to whatever transformation of features is most beneficial for subsequently deriving the object identity.

This transformation of features is not easy for a human to conceptualise, as we will see.

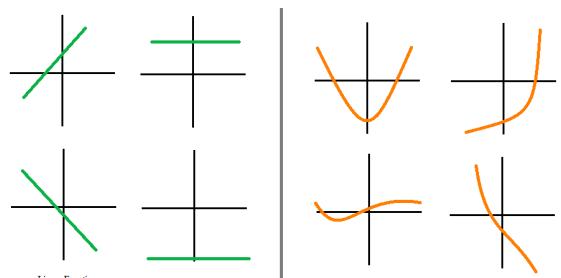
So, let's take a short break, and then look at how this is done.

TAKE A BREAK HERE:

The last part of this definition is ‘nonlinear processing units’.

But what does nonlinear mean, why is that necessary, and how is it achieved here?

Nonlinear functions



$$Y = A^*X + B$$

$$Y \neq A^*X + B$$

25. (Y is any other function of X)

In a linear function, the output (Y) of the function is simply the input (X) multiplied by a constant (A) and then added to another constant (B). The multiplier can be positive, negative or zero.

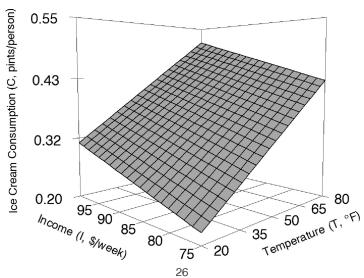
In a nonlinear function, there can be any other relationship between X and Y.

There must still be a relationship, Y is still a function of X, i.e. Y changes with X in some predictable way.

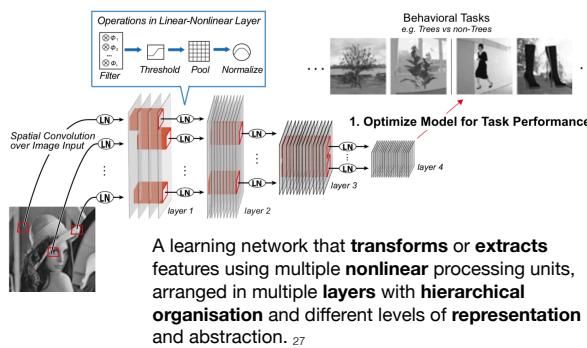
So we can see that non-linear functions can do a lot of things that linear functions can't.

Why nonlinear functions?

$$Y = B + A_1 * X_1 + A_2 * X_2 + \dots + A_p * X_p$$

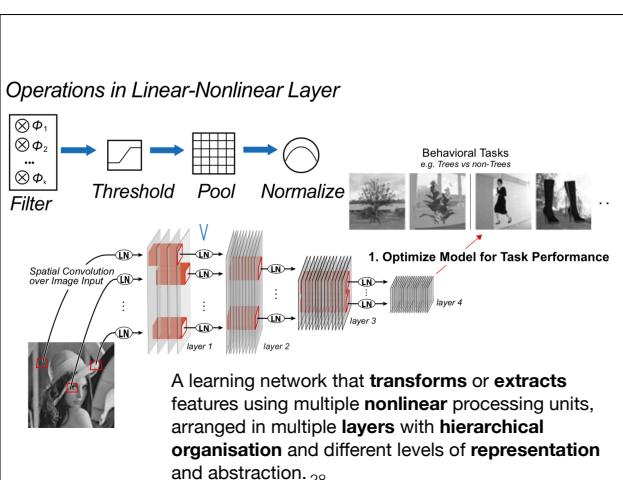


A deep network for object recognition



But in the context of deep networks, there is a more important problem with linear functions. The many layers of a deep network repeatedly perform functions on the output of each stage, joining these functions together to add effects of multiple features of the input. This becomes much harder to visualise, but this example, using only two input, gives an idea of the effect. Joining multiple linear functions (by addition or multiplication) always results in a linear function. That function can have multiple inputs, but the output is always a linear function of those inputs.

In this network, the inputs are the brightnesses of each image pixel. There is no way these can be multiplied and summed together to give the likelihood this is an image of a tree. Because, as we have seen, there is remarkably little relationship between an object's identity and the image it produces on the camera sensor. Indeed, any operation that can be done with only linear functions of the input can be straightforwardly described by formal mathematical rules, so is not a good use for deep networks.



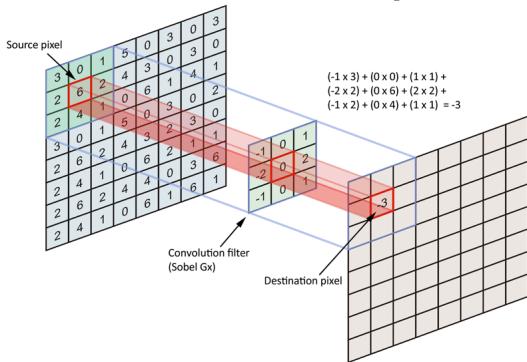
In our example, we have a complex nonlinear function with four operations or processing steps. These are filter, threshold, pool and normalise. These are very important to understand, so let's look at these steps in turn.

The output of one operation feeds into the next. This is true whether the operations are in the same layer or different layers. There is no step in this sequence of operations that has any special status, so the repeating sequence of these four operations effectively forms the layer. In some implementations, the result of each operation is seen as a layer. In your labs, filter and threshold are one command, while pooling is another, and normalisation is done throughout. Note that this is described as a linear-

nonlinear layer, which may be a little confusing.

The nonlinear function threshold is very important, but filter and normalise functions are linear. Pooling is optional

The filter/convolve operation



29

The filter or convolve operation is perhaps the most computationally important.

At the first processing layer, the input (or source) is a pixel map, a bitmap of the image giving the brightness of each pixel. In the simplest case, a grayscale image is used.

The convolution step looks for a pattern in a group of neighbouring pixels that corresponds to the convolution filter. For this filter, this would be dark on the left (low numbers) and light on the right (high numbers).

Using matrix multiplication, this filter is multiplied by a group of input pixels with a particular position, giving the match between the filter and a small part of the input image.

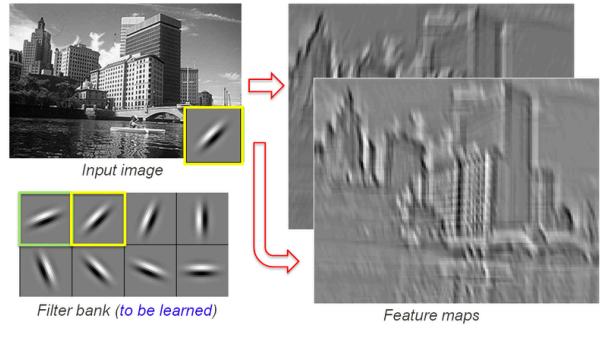
If the source pixels follow this filter pattern (light on the right, dark on the left), a high value will result. If the input area is all the same brightness, the result will be zero. If the source pixels are opposite to the filter (light on the right) the result will be negative.

Here, the match is poor: the pixel lightness on the left of the input source area is higher than on the right, so the output destination pixel gets assigned a low value.

Then, the filter is moved by one pixel in the input image to give a value for the next output pixel.

All source positions are multiplied by the filter to fill in all the destination pixels. A matrix multiplication function called 'convolution' does this efficiently, so this is often called a convolve/convolution step, leading to the term deep convolutional network.

The filter/convolve operation



The filter we just used is only one example to illustrate the principle.

In fact, a large set of filters are used in parallel to produce multiple maps of where the filter pattern is seen, often called feature maps.

Each 'pixel' shown in each of these feature maps is no longer a pixel in the image, it is an abstraction of the pattern across a group of pixels.

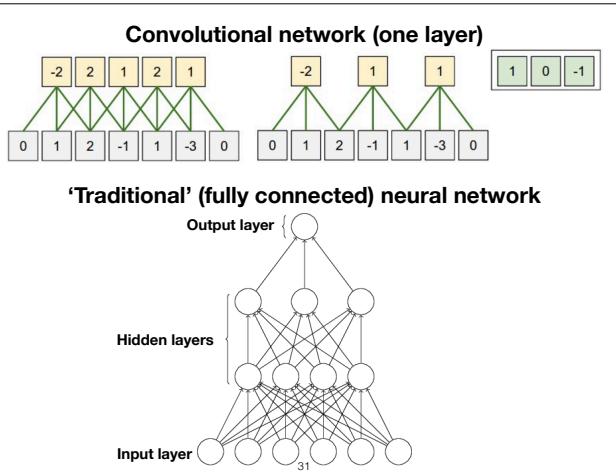
The pixel in the feature map represents the activity of a processing unit or artificial 'neuron'

Here we see a set of eight filters with larger extents and a range of different orientations.

These have been set manually in the first image analysis layer, because we know edge detectors with different orientations are an important early stage of computer and human vision systems.

The relevant filters can also be determined by machine learning, particularly for later layers where the best filter choice is not so obvious. We'll look at that possibility later.

The size of these filters will also determine the spatial scale of the feature we can detect. This is normally set manually using a small value for computational efficiency. Higher layers can detect larger scale features.



So, how is a convolutional network different from other types of neural network?

Here we see a 1-dimensional convolutional network. Here, the same filter is used at all positions in the network. Activation of each upper layer node depends on the product of this same filter convolved with a spatially limited set of lower layer nodes.

In a traditional neural net, each node in the upper layer is connected to all nodes in the previous layer, with no constraint on the spatial spread of links between them, and also no constraint on the pattern of links.

In a convolutional network, different patterns of filters are represented in different feature maps, but in a traditional neural network any 'filter' can effectively emerge for any node, so there is no need for multiple feature maps.

The constraint on the spatial spread of links between layers is particularly appropriate where the input layer nodes have a meaningful spatial relationship between them.

The most obvious examples of meaningful spatial relationships are in image processing, where the input layer units are image pixels and the feature maps each give the spatial distribution of the features described by the filters. We'll look at some other cases shortly.

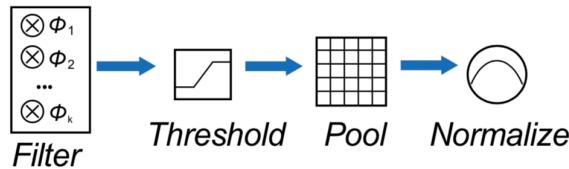
In the image example, the layers and filters are 2-dimensional because the input image is 2-dimensional. Here we are looking at 1-dimensional examples.

But in a fully connected network, every unit in one layer is connected to every unit in the next. Because of this, spatial relationship do not effect the network's activity.

A side effect of this is that the number of dimensions in the input is irrelevant in a fully connected network, so we can simply view each network layer as 1-dimensional.

The threshold/rectification operation

Operations in Linear-Nonlinear Layer



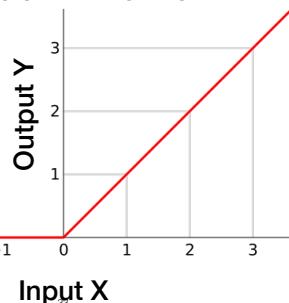
$$Y = f(X) = \max(0, X)$$

32

The threshold/rectification operation

an activation function using a rectified linear unit (ReLU)

$$Y = f(X) = \max(0, X)$$



33

The nonlinearity is introduced by the threshold or rectification operation.

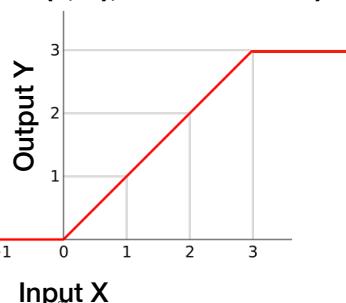
The goal of this operation is to only activate the output feature map if its value reaches a certain level, or threshold.

If we use filters that have a mean of zero, the threshold is typically zero.

The threshold/rectification operation

an activation function using a rectified linear unit (ReLU)

$$Y = f(X) = \min(\max(0, X), \text{MaxActivation})$$

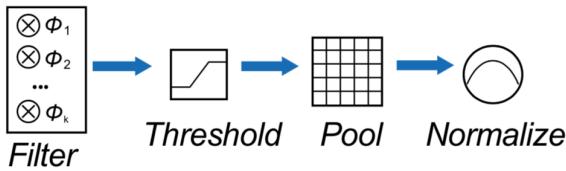


34

In practice, the maximum firing rate of a biological neuron is rarely reached, but a high maximum output is sometimes included, particularly if simulating biological neural systems.

The pooling operation

Operations in Linear-Nonlinear Layer



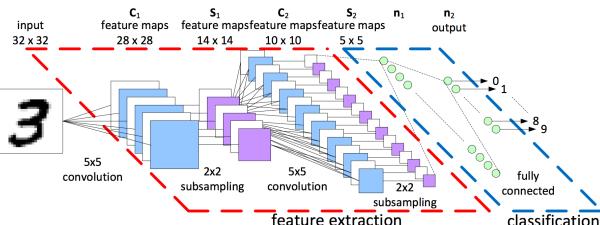
35

As a result of the filter operation, the response of each unit depends on several neighbouring inputs. So the units after filtering respond to a certain area of the input image, and the activation of neighbouring units will often be similar.

After several filter steps, each integrating inputs over an area, each unit will respond very similarly to an extensive area of the input. So neighbouring units are representing very similar information.

The pooling operation therefore downsamples the units to improve computational efficiency.

The pooling operation



36

Furthermore, from one input image, we have gone to several feature maps (C_1 here) at the filter/convolve operation.

The next filter/convolve operation will turn each of these into yet more feature maps. So to avoid an explosion of computational load, it's very important to reduce the size of these maps.

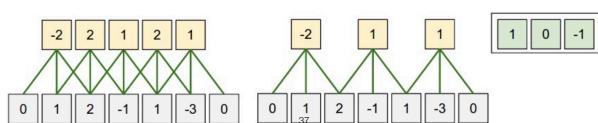
The pooling operation

Feature Map

6	4	8	5
5	4	5	8
3	6	7	7
7	9	7	2

Max-Pooling

6	8
9	7



This pooling operation is typically a simple max operation, taking the maximum of a square of 2×2 neighbouring units of the feature map.

This is very similar to the 'stride' parameter in the filter/convolve operation, so it is more efficient to increase the stride and skip the pooling.

BUT increasing stride may miss the maximum value, indeed both of the '2's are missed in this example.

The maximum value will be important for subsequent filtering.

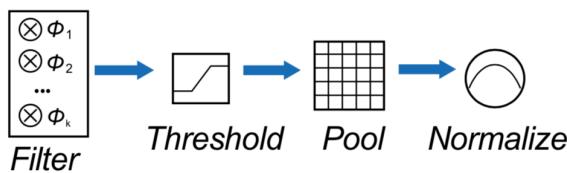
So the pooling operation discards some data in favour of computational efficiency.

As computers and deep network implementations become faster and more efficient, it should be less necessary to have

pooling layers. This would generally improve network performance but reduce speed.

The normalisation operation

Operations in Linear-Nonlinear Layer



38

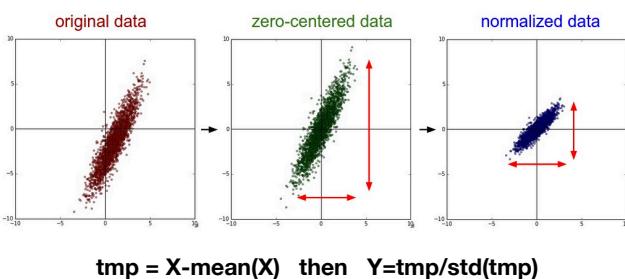
The threshold and pool operations use max functions.

As a result, even if the convolution filter has a mean of zero, by the pool stage we have a mean activation above zero and an arbitrary range.

Furthermore, this range be very different between feature maps, effectively weighting some feature maps to contribute more to the result than others.

Subsequent layers will have a problem here because subsequent filtering steps operate across multiple feature maps that may contribute very differently.

The normalisation operation



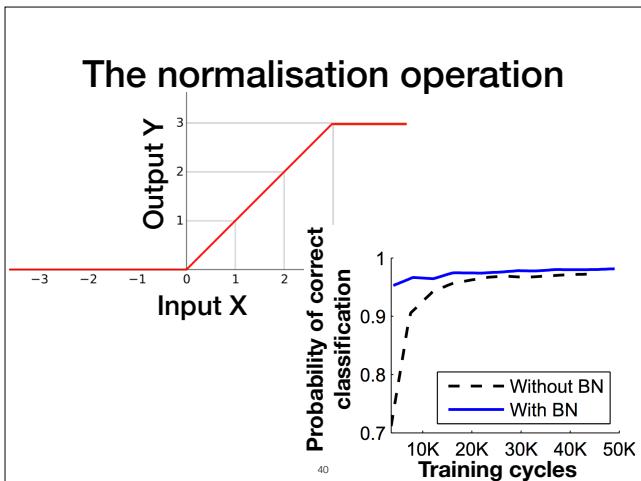
39

So the normalisation operation linearly scales the data to have a mean of zero activation for each feature map's responses to all images.

The first step of normalisation is the subtract the mean response of each feature map from all responses, i.e. to zero-center the data.

The next step is to divide the result by its standard deviation.

This makes the normalised data have a mean of zero and a standard deviation of one.



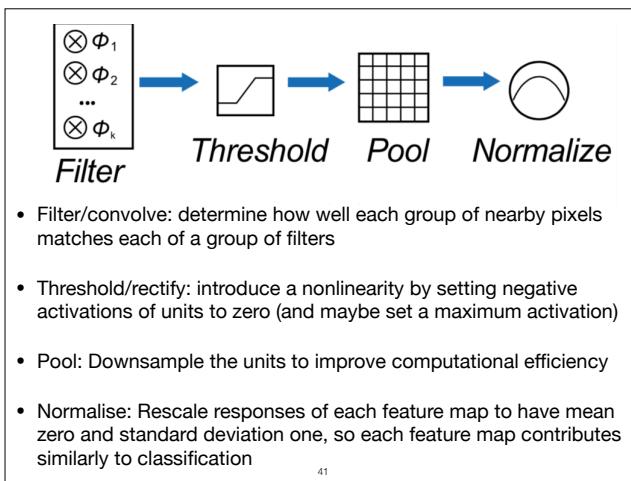
This is necessary for both theoretical and practical reasons

First, machine learning generally assumes that data reflects measurements of independent and identically-distributed (IID) variables. Normalisation forces identical distributions.

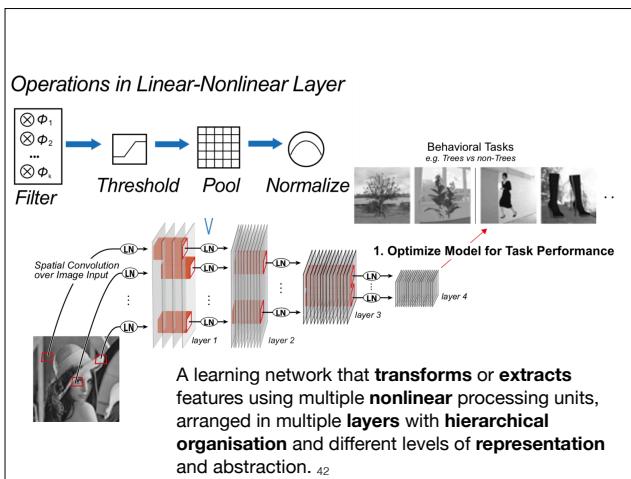
Second, if the activation function depends whether the unit's response is above or below zero, having zero-mean inputs and zero-mean filters, about half of the units will be active and half inactive. This even split of activation is a very efficient way to store information in a network of limited size.

Third, having the same range for all feature maps and layers means the same maximum threshold in the activation function can be used throughout the network.

Fourth, as a result of these considerations and other technical considerations, training rates are far better after normalisation, and final classification accuracy



So, we have now done one layer of deep network operations on an image. Let's summarise.

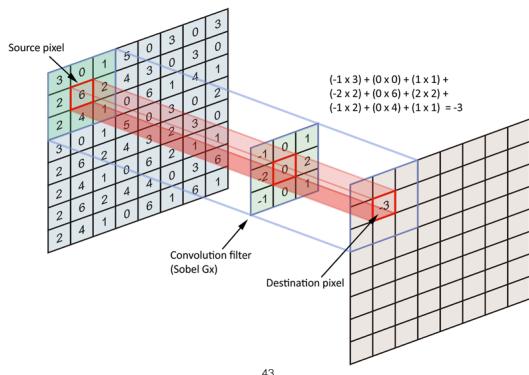


In a deep network, there are several layers performing similar operations and transformations of features, which all follow the same principles.

Subsequent layers generally use the same operations in the same way, but the filter/convolve operation differs between the first layer and subsequent layers.

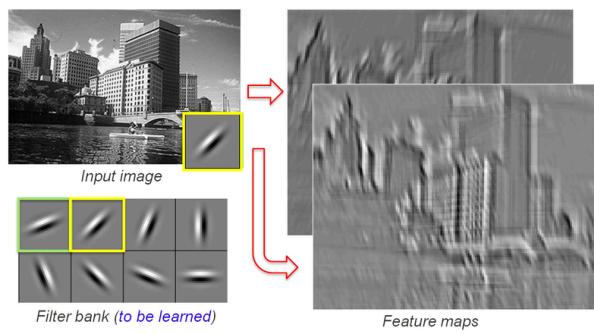
So far, we have only looked at the first layer, the simplest.

The filter/convolve operation (again)



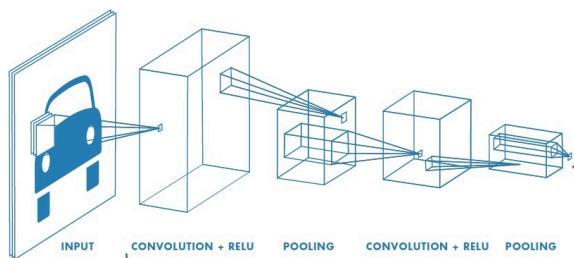
In the first layer, we are convolving a 2-dimensional input image with a two-dimensional filter to give a 2-dimensional feature map.

The filter/convolve operation (again)



But we do this for several filters, so the next layer is three-dimensional, with the third dimension corresponding to the filters used.

The filter/convolve operation (again)



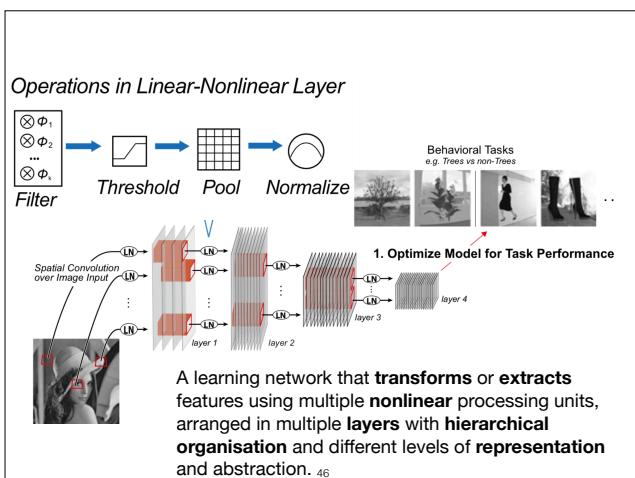
So in subsequent filter operations, the filter is also three-dimensional. It will normally span ALL the feature maps in its input to give the response at a single point in a single feature map of the next layer.

Remember that the result of a matrix multiplication can be a single number regardless of the size of the inputs to the multiplication.

This is even true for some input images: colour images are also treated as three dimensional, with the third dimension being the three colour channels.

So for colour images, the filters used are also three dimensional, with potentially different spatial responses required from each colour. So we can see the input image as a special case of the inputs to all layers. Any

convolution step generally considers filter patterns that cross all of its input feature maps.



As we get higher up the network, these filters get harder to understand in two important ways. First, the filter shape crosses multiple independent feature maps. An edge detector applied to an image is easy enough to conceptualise, but such a high-dimensional filter is harder to conceptualise.

Second, the input feature maps become more abstract. It gets very hard to conceptualise what feature is represented.

So both the filter and the feature map become a pattern within a pattern within a pattern. It's very hard to conceptualise these abstract, higher order patterns. Conveniently, we don't need to: the computer does this for us.

Inverting an object recognition DCNN



<http://thepsychreport.com/technology-2/googles-psychedelic-art-this-is-your-computer-brain-on-drugs/>

But it is important to understand that deep networks have increasingly complex representations of objects in the images, over several network layers.

To get a feel for what features are represented in different at different levels, we can ‘invert’ the network, activating different feature maps and their connections back to the original points in the image.

This imposed the features that each layer detected onto the original image. Here they use the example image of George Seurat's *Sunday on La Grande Jatte*. Successive layers find increasingly complex relationships between points in the original image, following correlations and patterns found in natural images.

Shared weights

- Filters generally have a single set of weights for all positions in the feature map because:
 - If a feature is useful to compute at one position, it is probably also useful at another position
 - The filter values are weights that need to be learned. It is very computationally demanding to do this if the set is too large.
 - The convolution operation is a very fast matrix function. If filters are not fixed, the convolution operation cannot be used

48

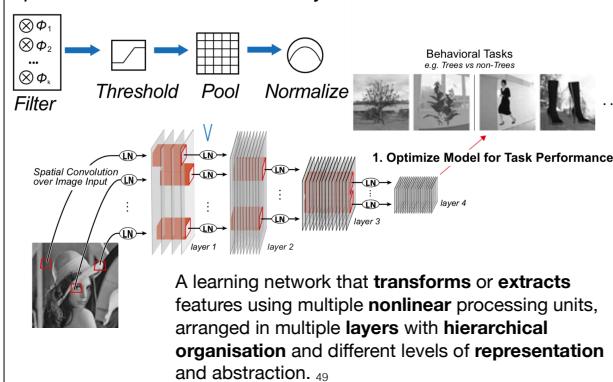
One filter is convolved with the previous feature map stack to give one new feature map.

It would also be possible for the filter to change at different positions in the feature map.

However, in artificial deep networks, generally a single filter is used across all positions, for three very good reasons.

AT END: We will see that these constraints do not apply in biological deep networks.

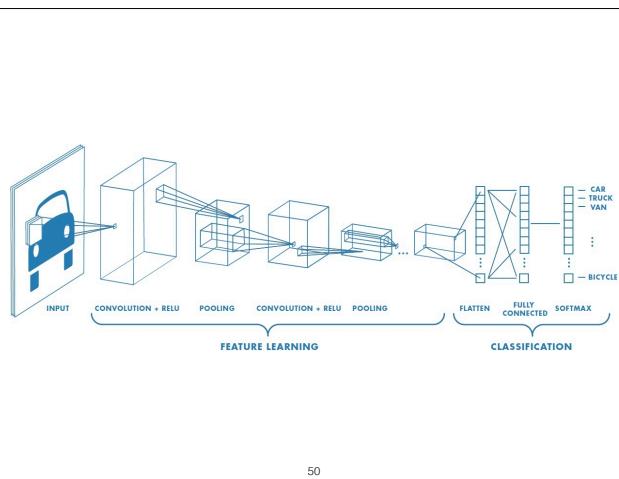
Operations in Linear-Nonlinear Layer



As we get higher in the network, more spatial integration occurs as a result of repeated convolution and pooling, so the spatial dimensions of the feature maps shrink.

At the same time, the number of interesting feature combinations increases, so the feature maps become increasingly narrow, but stacked increasingly high.

In the end, some classification or decision must be made, which is the fundamental goal of the network.



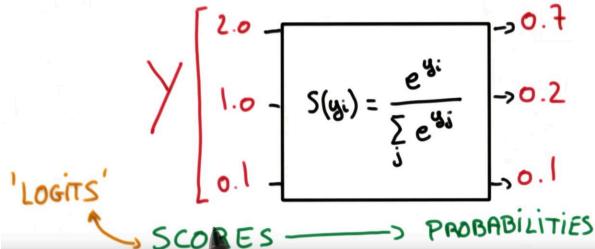
The last stage of the network, after several convolution layers, uses the activity in the final convolutional layer for classification. Here, the spatial relationships are first discarded, ‘flattening’ the last feature map into a line of independent units.

Each of these is connected to all the others, giving a fully-connected layer that looks at all links between every unit and all the others.

The activation pattern in these top-layer units is then associated with a particular output, a classification or label that describes the input image. Essentially the labels see which pattern they were trained on resembles the top-layer response pattern most closely.

The softmax operation

SOFTMAX



51

So, the weights through our network will transform each input image into a ‘score’, reflecting the match between the top layer’s pattern of activation by previous examples of each category.

This score must then be converted to a probability that this input image falls into each category.

This is almost always done with the ‘softmax’ function, or normalised exponential function. That is, constant e (the natural log base) raised to the power of the score, divided by the sum of these exponents for all scores to normalise the probabilities to sum up to one. (The math is not particularly important to know)

So far, we have carefully ignored what filters are used in higher levels. It is straightforward for the human researcher to design a simple filter, like an edge detector, to operate on an early layer, like the input image. But when we get to more abstract filters operating over multiple feature maps whose responses are already hard to conceptualise, we can no longer design a filter.

Instead, the filter structures are targets for machine learning. Indeed, the convolution filters are the main links between different layers of our network, so they effectively form the weights of connections between the nodes in a neural network.

Here, the nodes are pixels in a feature map, and the connections between these are filters.

So, to learn the weight of connections, the network learns the structure of the filters.

We constrain the filter’s size, maybe 3x3 or 5x5 pixels spanning all feature maps in a layer.

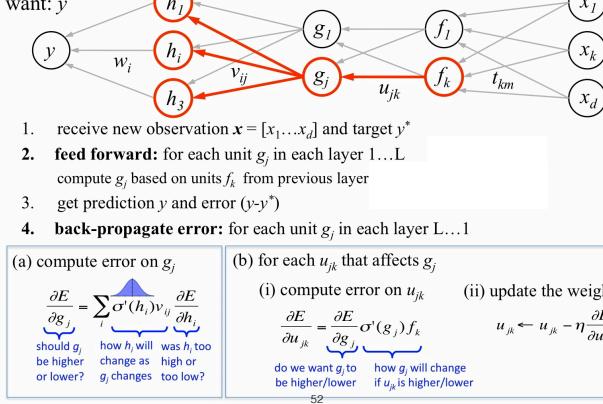
But the filter’s structure, the weights from the different pixels across the feature maps, is learned.

As we have seen, the filter is a limited group of connections, so the learning is much simpler than in a fully-connected network.

But just like learning in a fully-connected neural net, this uses backpropagation of error, a mathematically complex process that is notoriously difficult to understand and explain.

I will simplify as much as possible. Again, for this class, the math is not as important as what it is achieving.

Back-propagation



Deep learning in artificial neural networks

- Useful for achieving tasks that are difficult to describe formally
 - Difficult for computers, intuitive for humans
- A form of machine learning performing multiple sequential nonlinear feature transformations in hierarchical layers
- Between each feature map layer, a few simple operations
 - Convolution checks each position’s match with a specific filter kernel
 - Thresholding introduces a nonlinearity
 - Pooling downsamples the image, taking the maximum
 - Normalisation rescales responses so each feature map contributes similarly
- Final fully-connected layer links pattern of most abstracted, top-level features to required response
- Softmax determines probability of desired response
- Match or conflict of expected and actual outputs used as the basis for backpropagation of error (complex mathematical process)
 - Adjusts filter structure: link between layers, machine learning target

Before our next class, please think for yourself what the biological correlates of the filter/convolve, threshold, pooling and normalisation might be.

Also think about the full-connected last layer, the final classification, and learning by backpropagation.

Think whether these are feasible in biological neural networks, and whether they are necessary at all.

Lab Assignment

- In groups of 4, students will complete two lab assignments to build simple deep learning systems to solve computer vision and linguistics tasks. Students will be graded on two written reports of their work. Grades depend on depth and completion (20% each).
- For each set of questions:
 - First, work on the questions alone
 - Ask group members for help if you get stuck
 - Compare and discuss answers to agree a common single answer
 - Add this to a shared, combined final answers document (Google docs works well)
 - Show this to your teacher and discuss. You will be graded here.
 - Update answers if needed, show changes to teacher next time you speak.

54

Lab Assignment

- Intended to provide social contact and interaction
 - Use video where bandwidth allows, but this can slow your computer down a lot
- Each group has a single assigned teacher
 - Communicate in your Teams channel, because your teacher is also on this channel.
 - There is no need to make a WhatsApp group, Teams does the same things and more
 - We expect you to work with your group and teacher, in class time
 - Turn your video on when talking with your teacher, if possible
- This is a collaboration
 - The whole group usually gets the same grade
 - Let your teacher know if one of your group isn't participating
- There is also an individual assignment related to each group assignment
 - To distinguish the members of the group

55