

Artificial Intelligence / Machine Learning

One of the fathers of AI is worried about its future

Yoshua Bengio wants to stop talk of an AI arms race and make the technology more accessible to the developing world.

by Will Knight Nov 17, 2018

Yoshua Bengio is a grand master of modern artificial intelligence.

Alongside Geoff Hinton and Yann LeCun, Bengio is famous for championing a technique known as deep learning that in recent years has gone from an academic curiosity to one of the most powerful technologies on the planet.

Deep learning involves feeding data to large neural networks that crudely simulate the human brain, and it has proved incredibly powerful and effective for all sorts of practical tasks, from voice recognition and image classification to controlling self-driving cars and automating business decisions.

Bengio has resisted the lure of any big tech company. While Hinton and LeCun joined Google and Facebook, respectively, he remains a full-time professor at the University of Montreal. (He did, however, cofound Element AI in 2016, and it has built a very successful business helping big companies explore the commercial applications of AI research.)

Bengio met with MIT Technology Review's senior editor for AI, Will Knight, at an MIT event recently.

What do you make of the idea that there's an AI race between different countries?

I don't like it. I don't think it's the right way to do it.

MIT Technology Review

- What do you make of the idea that there's an AI race between different countries?
- Are there ways to foster more collaboration between countries?
- Are you worried about just a few AI companies, in the West and perhaps China, dominating the field of AI?
- There has been a lot of controversy over military uses of AI. Where do you stand on that?
- Even non-lethal uses of AI?
- Shouldn't AI experts work with the military to ensure this happens?
- What are you most excited about in terms of new AI research?
- You mention causality—in other words, grasping not just patterns in data but why something happens. Why is that important, and why is it so hard?

1

Google DeepMind



AlphaZero taught itself to play three different games

AlphaGo Vs Lee Sedol



AlphaZero

Game	Opponent	W: Wins	D: Draws	L: Losses
Chess	AlphaZero vs. Stockfish	29.0%	70.6%	0.4%
Shogi	AlphaZero vs. Elmo	84.2%	2.2%	13.6%
Go	AlphaZero vs. AGO	68.9%	31.1%	0%

[Source] <https://deepmind.com/blog/alphazero-shedding-new-light-grand-games-chess-shogi-and-go/>

Birkbeck, University of London

2

© Copyright 2019

2

BBC Sign in News Sport Weather iPlayer Sounds |

NEWS

Home UK World Business Election 2019 Tech Science Health Family & Education Technology

Go master quits because AI 'cannot be defeated'

27 November 2019 f t e Share



GETTY IMAGES

A master player of the Chinese strategy game Go has decided to retire, due to the rise of artificial intelligence that "cannot be defeated".

Lee Se-dol is the only human to ever beat the AlphaGo software developed by Google's sister company Deepmind.

In 2016, he took part in a five-match showdown against AlphaGo, losing four times but beating the computer once.

The South Korean said he had decided to retire after realising: "I'm not at the top even if I become the number one."

"There is an entity that cannot be defeated," the 18-time world Go champion told South Korea's Yonhap news agency.

Lee Se-dol is considered to be one of the greatest Go players of the modern era.

The 36-year-old former world champion started playing at the age of five, and turned pro just seven years later.

His defeat by the AlphaGo software was seen as a landmark moment for artificial intelligence.

3



Deep Learning I – Convolutional Neural Network

Dr Paul Yoo

Dept CSIS
28/11/19

Birkbeck, University of London 4 © Copyright 2019

4

Quiz



What is impersonation attack?

- A. an attack that allows an attacker to supply untrusted input to a program, which gets processed by an interpreter as part of a command or query which alters the course of execution of that program.
- B.** an attack in which an adversary successfully assumes the identity of one of the legitimate parties in a system or in a communications protocol.
- C. an attack that is designed to bring a network or service down by flooding it with large amounts of traffic.
- D. an attack in which a system is monitored and sometimes scanned for open ports and vulnerabilities.

Quiz



What are characteristics of network (wireless) based IDS?

- 1.** They look for attack signatures in network traffic
2. Filter decides which traffic will not be discarded or passed
3. It is programmed to interpret a certain series of packet
4. It models the normal usage of network as noise characterization

Quiz



What are the major components of the IDS?

- A. Detection engine
- B. Alert processing
- C. Report processing
- D. Packet feature extraction
- E. All of the above

Quiz



What are the reasons that ML is becoming a key weapon in the cyber security war?

- A. Because it is a popular technology
- B. Antivirus software is too expensive
- C. Model free properties
- D. Legacy antivirus technology is too slow to stop cyber-attack in time. Attacks are often weeks or months old by the time they are discovered.
- E. All of the above

Quiz



What are the ML algorithms used by D-FES and DETERed?

- A. Autoencoder
- B. RBFN
- C. Logistic Regression
- D. MLP
- E. All of the above

Quiz



Which of the below ML algorithms make the most money?

- A. Generative Adversarial Network
- B. Multi-layered Perceptron
- C. Logistic Regression
- D. Convolutional Neural Network
- E. Support Vector Machine
- F. Recurrent Neural Network
- G. Bayesian Network
- H. None of the above

Overview



We covered:

- Group Practical
- Application: IoT Intrusion Detection
- AWID (Kolias et al 2015)
- DEMISe (Parker et al 2019)
- Evaluation

We will cover:

- Convolutional Neural Network
 - Convolution in 2D and 3D
 - Non-linearity
 - Pooling
 - Classification
 - Other CNNs (Alexnet, VGG16)

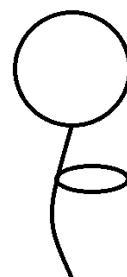
**Consider a boy who's learning drawing
and painting for the first time.**



He has a dad who is not very good at art.



After a while, he draws this,



Birkbeck, University of London

13

© Copyright 2019

13

Let me show you a few more pictures



This time he draws something like this.



Birkbeck, University of London

14

© Copyright 2019

14

His dad repeats the process several times



Slowly and steadily he learns many things. From the basic skeleton to its fine features. Like:

- Shades and Shapes of the petals
- Edges of the leaves
- Features of the stem
- Etc

A CNN has two main components



Convolution filters + a full connected network

Convolutional Filters

- the boy learnt about the following
 - the round thing on the top or the stem and the leaf shape – very basic
 - a collection of similar looking units or petals – basic
 - the variation in shades – a bit more complex
- CFs capture **low to high-level features** of the image.

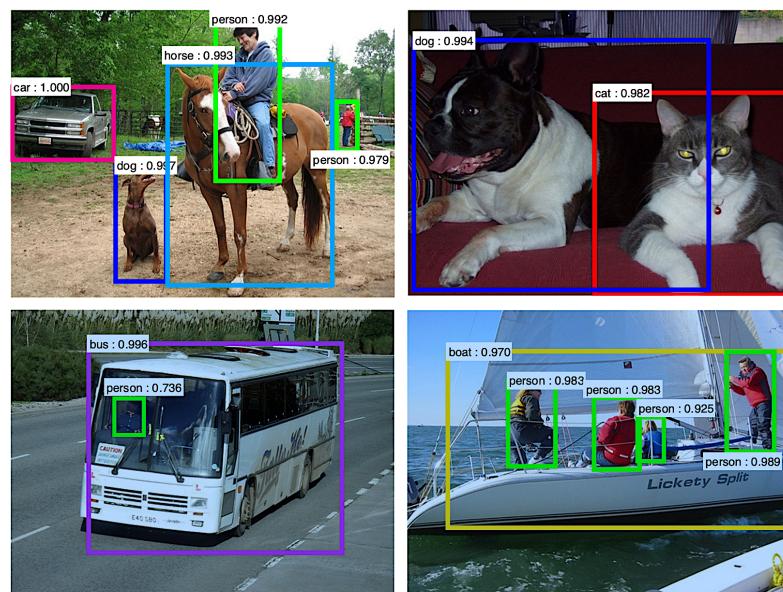
What are CNNs (a.k.a ConvNets) ?



CNNs are a category of neural networks that have proven very effective in areas such as image recognition and classification.

ConvNets

- have been successful in identifying faces (e.g. Facebook), objects and traffic signs apart from powering vision in robots and self-driving cars.
- is able to recognise scenes and the system is able to suggest relevant captions ("a football player is kicking a football ball")
- is used for recognising everyday objects, humans and animals.
- have been effective in several NLP tasks (such as sentence classification) as well.



Source: <https://arxiv.org/pdf/1506.01497v3.pdf>

PROC. OF THE IEEE, NOVEMBER 1998

1

Gradient-Based Learning Applied to Document Recognition

// LeNet was one of the very first CNNs which helped propel the field of Deep Learning.

Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner

Abstract—

Multilayer Neural Networks trained with the backpropagation algorithm constitute the best example of a successful Gradient-Based Learning technique. Given an appropriate network architecture, Gradient-Based Learning algorithms can be used to synthesize a complex decision surface that can classify high-dimensional patterns such as handwritten characters with minimal error. This paper reviews various methods applied to handwritten character recognition and compares them on a standard handwritten digit recognition task. Convolutional Neural Networks, that are specifically designed to deal with the variability of 2D shapes, are shown to outperform all other techniques.

Real-life document recognition systems are composed of multiple modules including field extraction, segmentation, recognition, and language modeling. A new learning paradigm, called Graph Transformer Networks (GTN), allows such multi-module systems to be trained globally using Gradient-Based methods so as to minimize an overall performance measure.

Two systems for on-line handwriting recognition are described. Experiments demonstrate the advantage of global training, and the flexibility of Graph Transformer Networks.

A Graph Transformer Network for reading bank check is also described. It uses Convolutional Neural Network character recognizers combined with global training techniques to provide record accuracy on business and personal checks. It is deployed commercially and reads several million checks per day.

Keywords— Neural Networks, OCR, Document Recognition, Machine Learning, Gradient-Based Learning, Convolutional Neural Networks, Graph Transformer Networks, Finite State Transducers.

Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. Proceedings of the IEEE, November 1998

19

The LeNet Architecture (1990s)

AT&T **LeNet 5** RESEARCH

answer: 1

1
1
1

// the LeNet architecture was used mainly for character recognition tasks such as reading zip codes, digits, etc.

// trained on grey-scale images

Birkbeck, University of London

20

© Copyright 2019

20

CNN

There are four main operations (Conv., Non linearity, Pooling and Classification) – building blocks of CNN!

Convolution + ReLU Pooling Convolution + ReLU Pooling Fully Connected Fully Connected Output Predictions

dog (0.01)
cat (0.04)
boat (0.94)
bird (0.02)

// correctly assigns the highest probability for boat (0.94) among all four categories.

/ this CNN architecture is similar in architecture to the original LeNet and classifies an input image into four categories: dog, cat, boat or bird

Birkbeck, University of London 21 © Copyright 2019

21

Images

Every image can be represented as a matrix of pixel values.

Channel

- An image from a standard digital camera will have **three channels** – red, green and blue.
- Three **2d-matrices stacked over each other** (one for each color), each having pixel values in the **range 0 to 255**.

A grayscale image has just one channel.

- A **single 2d matrix** representing an image.
- The value of each pixel in the matrix will range from 0 to 255 – zero indicating black and 255 indicating white.

Birkbeck, University of London 22 © Copyright 2019

22

Convolution

//
Python: conv.forward
TensorFlow: tf.nn.conv2d
Keras: Conv2D



The purpose of Convolution is to extract features from the input image.

// the convolution operation. the output matrix is called Convolved Feature or Feature Map.

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Image

// filter "sees" only a part of the input image in each stride.

4		

Convolved Feature

// consider a 5 x 5 image whose pixel values are only 0 and 1

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

// consider another 3 x 3 matrix // called 'filter' or 'kernel' or 'feature detector'

1. Slide the filter (orange matrix) over our original image (green) by 1 pixel (also called 'stride') and for every position

2. Compute element wise multiplication (between the two matrices)

3. Add the multiplication outputs to get the final integer which forms a single element of the output matrix (pink).

Birkbeck, University of London 23 © Copyright 2019

23

Filter



Different values of the filter matrix will produce different Feature Maps for the same input image.

Operation	Filter	Convolved Image
Identity	$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	
Edge detection	$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$	
	$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$	
	$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$	
Sharpen	$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$	
Box blur (normalized)	$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$	
Gaussian blur (approximation)	$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$	

// different filters can detect different features from an image, for example edges, curves etc.

Input

The Convolution Operation

// this captures the local dependencies in the original image.

Source: https://cs.nyu.edu/~fergus/tutorials/deep_learning_cvpr12/

Birkbeck, University of London 24 © Copyright 2019

24

Vertical & horizontal edges

The diagram shows a grayscale image of people and a bicycle. Two convolutional filters are applied to it:

- Vertical edges:** A filter with kernel $\begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}$. The resulting output image highlights vertical edges.
- Horizontal edges:** A filter with kernel $\begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$. The resulting output image highlights horizontal edges.

Birkbeck, University of London

25

© Copyright 2019

25

Convolution over volume

Multiple filters can be used to detect multiple features. # of output of the layer = # of filters in the layer.

// when the input has more than one channels (e.g. an RGB image), the filter should have matching number of channels.

Input: $6 \times 6 \times 3$ (3 channels)

Filter: $3 \times 3 \times 3$ (3 channels)

Result: $4 \times 4 \times 2$ (2 channels)

Parameters:
Size: $f=3$
#channels: $n_c=3$
Stride: $s=1$
Padding: $p=0$

<https://indoml.com>

Input: $6 \times 6 \times 3$

Filter 1: $3 \times 3 \times 3$ (3 channels)

Output: 4×4 (3 channels)

Filter 2: $3 \times 3 \times 3$ (3 channels)

Output: 4×4 (3 channels)

$(4 \times 4 \times 2) \times (3 \times 3 \times 3) = 864$

The total number of multiplications to calculate the result is $(4 \times 4) \times (3 \times 3 \times 3) = 432$

Birkbeck, University of London

26

© Copyright 2019

26

1 x 1 Convolution

This is convolution with 1×1 filter. The effect is to flatten or “merge” channels together, which can save computations later in the network:

Input

4	9	2	5	8	3
5	6	2	4	0	3
2	4	5	4	5	2
5	6	5	4	7	8
5	7	7	9	2	1
5	8	5	3	8	4

Filter

1	1	1	1	1
1	1	1	1	1
1	1	1	1	1
1	1	1	1	1
1	1	1	1	1

Result

4	9	2	5	8	3
5	6	2	4	0	3
2	4	5	4	5	2
5	6	5	4	7	8
5	7	7	9	2	1
5	8	5	3	8	4

Parameters:

- Size: $f = 1$
- #channels: $n_C = 5$
- Stride: $s = 1$

<https://indoml.com>

Birkbeck, University of London

27

© Copyright 2019

27

How to learn the values of filter?

- In practice, a CNN learns the values of these filters on its own during the training process
 - Specify parameters such as **# of filters**, **filter size**, **architecture of the network etc.** before the training process.
- The more **# of filters** we have, the more image features get extracted and the better our network becomes at recognising patterns in unseen images.

Birkbeck, University of London

28

© Copyright 2019

28

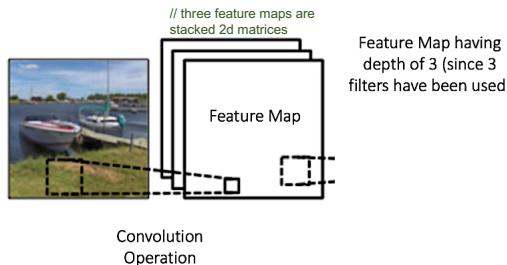


Feature Map

Decide three parameters the convolution step is performed (Depth, Stride and Padding)

Depth corresponds to the # of filters we use for the convolution operation.

- Perform convolution of the original boat image using **three distinct filters**, thus producing **three different feature maps** as shown.



Birkbeck, University of London

29

© Copyright 2019

29

Feature Map cont.

$$\begin{matrix} 2 & 3 & 7 & 4 & 6 & 2 & 5 \\ 6 & 9 & 8 & 7 & 4 & 3 \\ 3 & 4 & 8 & 3 & 8 & 9 & 7 \\ 7 & 8 & 3 & 6 & 6 & 3 & 4 \\ 4 & 2 & 1 & 8 & 3 & 4 & 6 \\ 3 & 2 & 4 & 1 & 9 & 8 & 3 \\ 0 & 1 & 3 & 9 & 2 & 1 & 4 \end{matrix} \quad \begin{bmatrix} 1 & 0 & 4 \\ 1 & 0 & 2 \\ -1 & 0 & 3 \end{bmatrix} = \left[\frac{n+2p-f}{s} + 1 \right] \times \left[\frac{n+2p-f}{s} + 1 \right]$$

*f x f filter
n x n image
stride = s
padding = p*

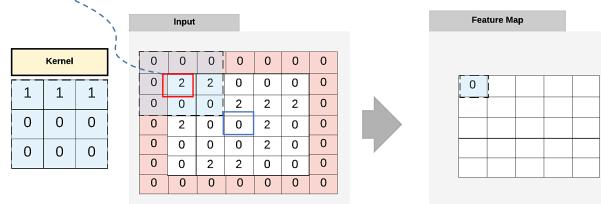


Stride is the number of pixels by which we slide our filter matrix over the input matrix.

- When the stride is 1 then we move the filters **one pixel at a time**.
- When the stride is 2, then the filters **jump 2 pixels at a time** as we slide them around.
- Having a **larger stride** will produce **smaller feature maps**.

Zero-padding.

- Problems with **"valid" convolution** (no padding) $\rightarrow n-f+1 \times n-f+1$
 - Shrinking output.
 - Throwing away info from edge – the corner pixels used only one of the outputs.
- **"Same" Convolution** : Pad so that output size is the same as input size.
 - Pad the input with zeros around the border, so that we can **apply the filter to bordering elements** of our input image matrix.
 - **Control the size of the feature maps.** $\rightarrow p = (f-1)/2$



Birkbeck, University of London

30

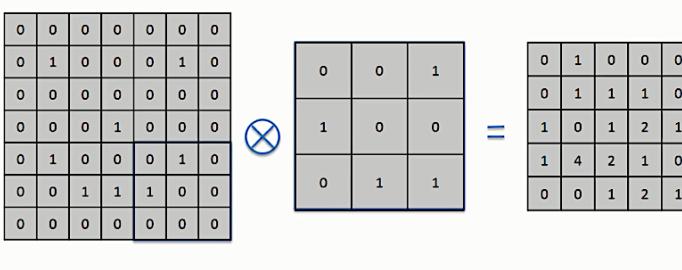
© Copyright 2019

30



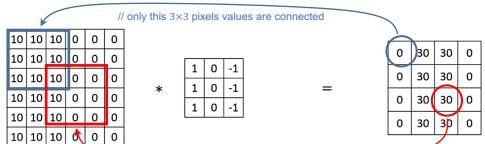
- What is the value of stride? $\left\lceil \frac{n+2p-f}{s} + 1 \right\rceil \times \left\lceil \frac{n+2p-f}{s} + 1 \right\rceil = \left\lceil \frac{7-3}{7} + 1 \right\rceil = 5$
- What are other names of feature detector?

Filter, kernel



Input Image Feature Detector Feature Map

Quiz



Why do we use the convolution operation?

- To extract features from the input image.
- Reduce # of parameters
 - Parameter sharing** – feature detector that's useful in one part of the image is probably useful in another part of the image.
 - Sparsity of connections** – in each layer, each output value depends only on a small number of inputs (less prone to overfitting)
- Reducing the size of the input image, the larger your strides, the smaller your feature map.
- Real images tend to be substantially larger and more complex than the 7 x 7 image – widen your strides.
- Easier to read.

Quiz



Do we lose information when using a feature detector?

- Yes – the feature map that we [end up with has fewer cells](#) and therefore [less information](#) than the original input image.
- However, the very purpose of the feature detector is to [sift through the information](#) in the input image and filter the parts that are integral to it and exclude the rest.
- It is meant to separate the wheat from the chaff.

Quiz



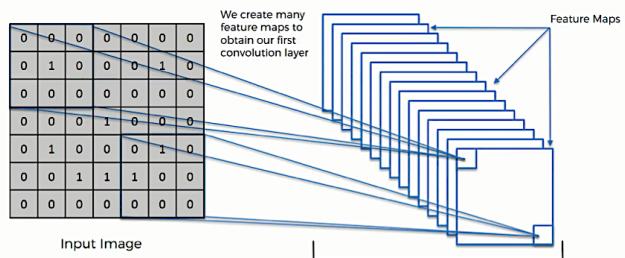
Why do we aim to reduce the input image to its essential features?

- What the convolution does is detect certain features
 - e.g. eyes and nose – you immediately know who you are looking at.
- These are all your brain needs to see in order to make its conclusion.
- CNNs operate in the similar way as your brain



Quiz

How to CNNs actually perform this operation?



- Through **training**, the network determines **what features it finds important** in order for it to be able to scan images and categorise them more accurately.
 - Based on that, it develops its feature detectors.
 - In many cases, the features considered by the network will be **unnoticeable to the human eye**, which is exactly why CNNs are so amazingly useful.
- // in reality, CNNs develop multiple feature detectors and use them to develop several feature maps which are referred to as convolutional layers.



Quiz

What are other uses of convolution matrices?

- The word "filter" is used in the same sense we use it when talking about [Instagram filters](#).
- You can actually use a convolution matrix [to adjust an image](#).

Sharpen:

$$\begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 1 & 5 & 1 & 0 & 0 \\ 0 & 1 & 5 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$



Blur:

$$\begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$



Edge Detect:

$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$



Non Linearity (ReLU)

Break up the linearity in the image.

- e.g. the transition between pixels, the borders, the colours
- ReLU has been used after every convolution operation.
- ReLU is an element wise operation (applied per pixel) and replaces all negative pixel values in the feature map by zero.

// convolution is a linear operation
– element wise matrix multiplication and addition

// tanh or sigmoid can also be used

Input Feature Map Rectified Feature Map

Black = negative; white = positive values

Only non-negative values

ReLU

// remove all the black elements from it, keeping only those carrying a positive value (the grey and white colors).

Birkbeck, University of London 37 © Copyright 2019

37

Pooling

Enables the CNNs to detect the cheetah when presented with the image in any manner.

Birkbeck, University of London 38 © Copyright 2019

38



Pooling cont.

To teach your CNN to recognise that despite all of these differences, the network needs to acquire a property that is known as “spatial variance.”

// the capability of detecting the object in the image without being confused by the differences in the image's textures, the distances from where they are shot or their angles.

// each posing differently in different settings and from different angles.



Image Source: Wikipedia

Birkbeck, University of London

39

© Copyright 2019

39



Pooling cont.

Spatial Pooling (a. k. a. subsampling or downsampling) reduces the dimensionality of each feature map but retains the most important information.

- Spatial Pooling can be of different types: Max, Average, Sum etc.
- In case of Max Pooling, we define a spatial neighborhood (for example, a 2×2 window) and take the largest element from the rectified feature map within that window.
- In case of Average Pooling, instead of taking the largest element we could also take the average of all elements in that window.
- In practice, Max Pooling has been shown to work better.

Birkbeck, University of London

40

© Copyright 2019

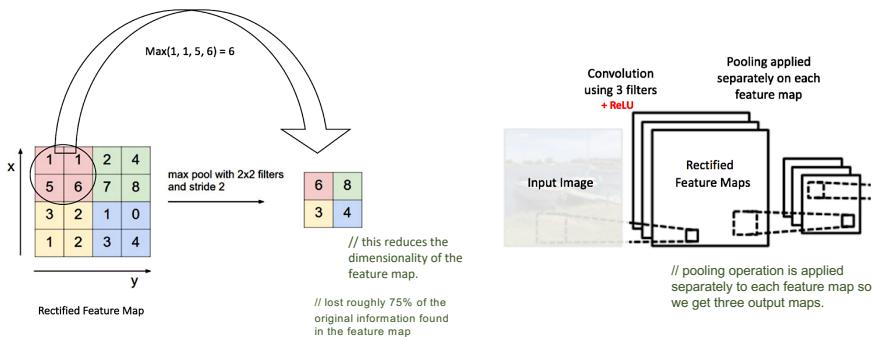
40

Pooling cont.



// generally, strides of two are most commonly used.

Slide 2×2 window by 2 cells (also called ‘stride’) and take the maximum value in each region. Two hyperparameters (f and s).



Birkbeck, University of London

41

© Copyright 2019

41

Pooling cont.



- makes the input representations (feature dimension) **smaller** and more **manageable**.
- **reduces the # of parameters and computations** in the network, therefore, controlling **overfitting**.
- makes the network **invariant to small transformations, distortions and translations** in the input image (a small distortion in input will not change the output of Pooling – since we take the maximum / average value in a local neighborhood).
- helps arrive at an almost scale invariant representation of our image (the exact term is “**equi-variant**”) – this is very powerful since we can **detect objects in an image no matter where they are located**.

Birkbeck, University of London

42

© Copyright 2019

42

Quiz



What does the pooling process provide CNN?

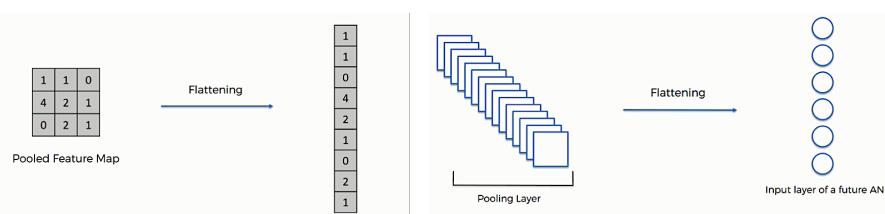
- Equivariant capability
- Speed up computation
- Reduce overfitting (abstract features could have less rules)

Flattening



We're going to need to insert this pooled feature map data into an ANN now.

- Flatten our pooled feature map into a column like in the image below.



// after the flattening step is that you end up with a long vector of input data that you then pass through the artificial neural network to have it processed further.

Three basic building blocks of CNN

The diagram illustrates the architecture of a Convolutional Neural Network (CNN). It starts with an input image of a boat on water. This image is processed by a sequence of layers:

- Convolution + ReLU**: The input image is processed by a convolution layer using six filters, resulting in six feature maps. A ReLU activation function is applied to each map.
- Pooling**: The feature maps are pooled to produce a total of six feature maps.
- Convolution + ReLU**: The output of the first pooling layer is processed by another convolution layer using six filters, resulting in six feature maps. A ReLU activation function is applied to each map.
- Pooling**: The feature maps are pooled to produce a total of six feature maps.
- Fully Connected**: The output of the second pooling layer is flattened and passed through a fully connected layer.
- Fully Connected**: The output of the fully connected layer is passed through another fully connected layer.
- Output Predictions**: The final output consists of four predicted classes with their respective probabilities: dog (0.01), cat (0.04), boat (0.94), and bird (0.02).

Annotations provide additional details:

- // performs Convolution on the output of the first Pooling Layer using six filters to produce a total of six feature maps.
- // ReLU is then applied individually on all of these six feature maps.
- // perform Max Pooling operation separately on each of the six rectified feature maps.
- Together these layers
 - extract the useful features from the images,
 - introduce non-linearity in the network and
 - reduce feature dimension
 - make the features somewhat equivariant to scale and translation
- The output of the 2nd Pooling Layer acts as an input to the Fully Connected Layer.

Birkbeck, University of London 45 © Copyright 2019

Fully Connected Layer

The Fully Connected Layer (FCL) is a traditional Multi-Layer Perceptron (MLP) that uses a softmax activation function in the output layer.

- The output from the convolutional and pooling layers represent **high-level features** of the input image.
- MLP uses these features for classifying the input image into various classes based on the training dataset.
- The sum of output probabilities from the Fully Connected Layer is 1.

Birkbeck, University of London 46 © Copyright 2019

Deep Learning using Linear Support Vector Machines

Yichuan Tang TANG@CS.TORONTO.EDU
Department of Computer Science, University of Toronto. Toronto, Ontario, Canada.

Abstract

Recently, fully-connected and convolutional neural networks have been trained to achieve state-of-the-art performance on a wide variety of tasks such as speech recognition, image classification, natural language processing, and bioinformatics. For classification tasks, most of these “deep learning” models employ the softmax activation function for prediction and minimize cross-entropy loss. In this paper, we demonstrate a small but consistent advantage of replacing the softmax layer with a linear support vector machine. Learning minimizes a margin-based loss instead of the cross-entropy loss. While there have been various combinations of neural nets and SVMs in prior art, our results using L2-SVMs show that by simply replacing softmax with linear SVMs gives significant gains on popular deep learning datasets MNIST, CIFAR-10, and the ICML 2013 Representation Learning Workshop’s face expression recognition challenge.

multistage process. In particular, a deep convolutional net is first trained using supervised/unsupervised objectives to learn good invariant latent representations. The corresponding hidden variables of data samples are then treated as input and fed into linear (or kernel) SVMs (Huang & LeCun, 2006; Lee et al., 2009; Quoc et al., 2010; Coates et al., 2011). This technique usually improves performance but the drawback is that lower level features are not been fine-tuned w.r.t. the SVM’s objective.

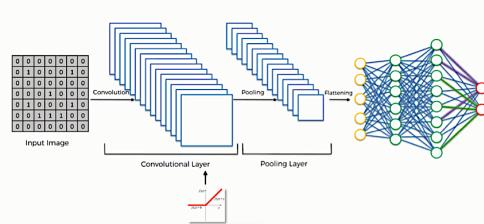
Other papers have also proposed similar models but with joint training of weights at lower layers using both standard neural nets as well as convolutional neural nets (Zhong & Ghosh, 2000; Collobert & Bengio, 2004; Nagi et al., 2012). In other related works, Weston et al. (2008) proposed a semi-supervised embedding algorithm for deep learning where the hinge loss is combined with the “contrastive loss” from siamese networks (Hadsell et al., 2006). Lower layer weights are learned using stochastic gradient descent. Vinyals et al. (2012) learns a recursive representation using linear SVMs at every layer, but without joint fine-tuning of the hidden representation.

47

A Quick Recap



- We start off with an **input image**.
- We apply **filters** (or feature maps) to the image, which gives us a convolutional layer.
- We then **break up the linearity** of that image using the rectifier function.
- The image becomes ready for the **pooling** step, the purpose of which is providing our CNN with the faculty of “spatial invariance”.
- After we’re done with pooling, we end up with a pooled feature map.
- We then **flatten** our pooled feature map before inserting into an ANN.
- Throughout this entire process, the **weights and the feature maps are trained** and repeatedly altered in order for the network to reach the optimal performance.



Birkbeck, University of London 48 © Copyright 2019

48

Putting it all together



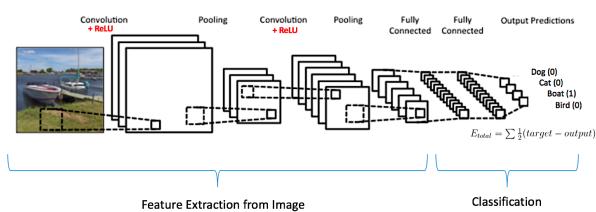
Training using backpropagation

- Convolution + Pooling layers act as Feature Extractors from the input image while Fully Connected layer acts as a classifier.

- For example,

Input Image = Boat

Target Vector = [0, 0, 1, 0]



Birkbeck, University of London

49

© Copyright 2019

49

- initialise all filters and parameters / weights with random values
- the network takes a training image as input, goes through the forward propagation step (convolution, ReLU and pooling operations along with forward propagation in the FC layer) and finds the output probabilities for each class.
 - lets say the output probabilities for the boat image above are [0.2, 0.4, 0.1, 0.3]
 - since **weights are randomly assigned** for the first training example, output probabilities are also random.
- calculate the total error at the output layer (summation over all 4 classes) – **cross entropy** or a mean squared error
 - Total Error = $\sum \frac{1}{2} (\text{target probability} - \text{output probability})^2$



Birkbeck, University of London

50

© Copyright 2019

50

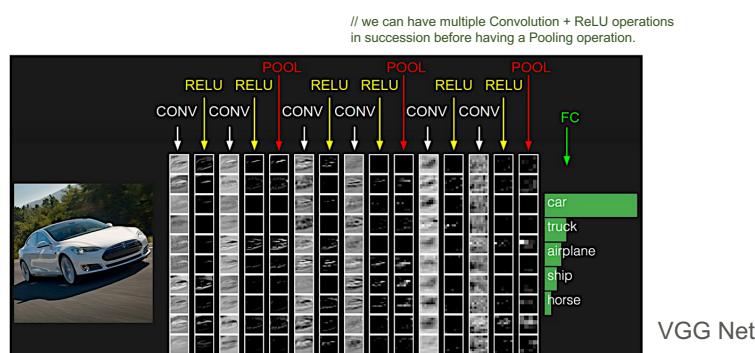
4. use **backpropagation** to calculate the *gradients* of the error with respect to all weights in the network and use **gradient descent** to update **all filter values / weights** and **parameter values** to minimise the output error.
- the weights are adjusted** in proportion to their contribution to the total error.
 - when the same image is input again, output probabilities might now be [0.1, 0.1, 0.7, 0.1], which is closer to the target vector [0, 0, 1, 0].
 - this means that the network has *learnt* to classify this particular image correctly by adjusting its weights / filters such that **the output error is reduced**.
 - parameters (# of filters, filter sizes, architecture of the network etc) have all been fixed before Step 1** and do not change during training process – only the values of the filter matrix and connection weights get updated.
5. repeat steps 2–4 with all images in the training set.



Convolution Design

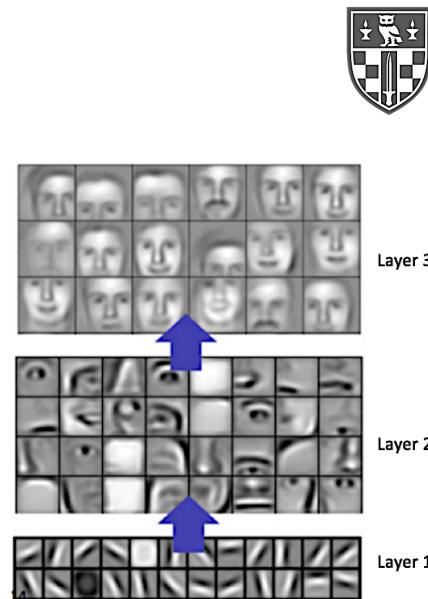


- In fact, some of the best performing CNNs today have **tens of Convolution and Pooling layers!**
- It is **not** necessary to have a Pooling layer after every Convolutional Layer.



Convolution Steps

- In general, the **more convolution steps** we have, the **more complicated features** our network will be able to learn to recognize.
- For example,
 - 1st layer learns to detect **edges** from raw pixels,
 - 2nd layer uses the edges to detect **simple shapes**
 - 3rd layer uses these shapes to detect **higher-level features** such as facial shapes.



Birkbeck, University of London

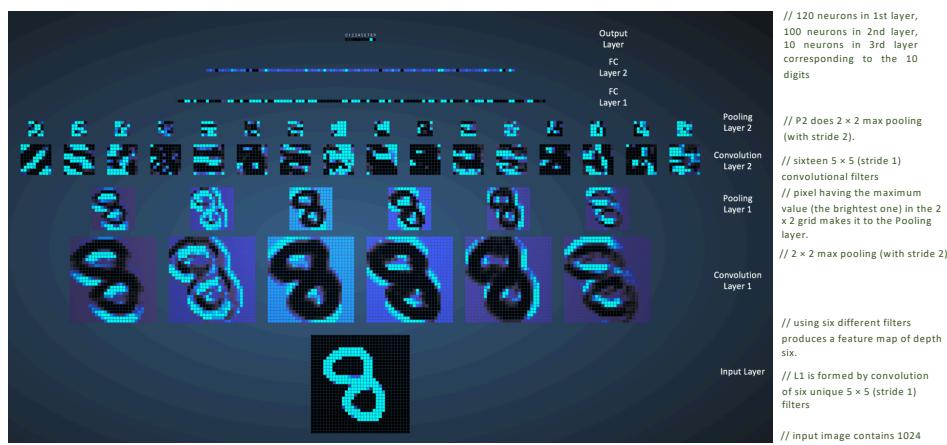
53

© Copyright 2019

53

Visualisations of a CNN trained on the MNIST

Harley, A.W., 2015, December. An interactive node-link visualization of convolutional neural networks. In *International Symposium on Visual Computing* (pp. 867-877). Springer, Cham.



Birkbeck, University of London

54

© Copyright 2019

54



Other CNNs

CNNs have been around since early 1990s.

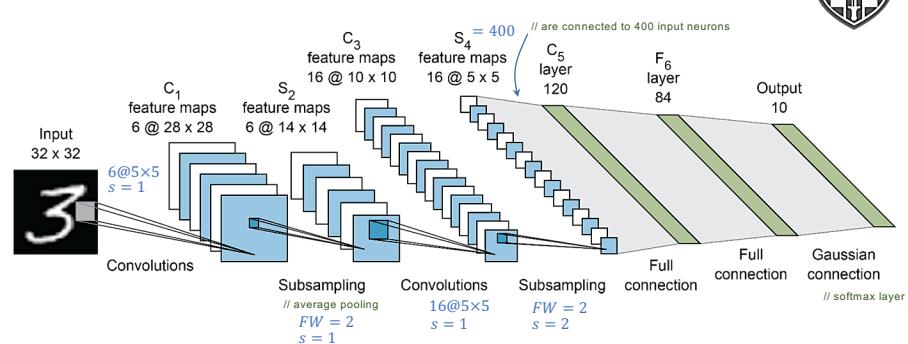
LeNet-5 (1990s)

- This pioneering work by Yann LeCun was named LeNet5 after many previous successful iterations since the year 1988!
 - There was no GPU to help training, and even CPUs were slow.
 - CNNs use sequence of 3 layers: convolution, pooling, non-linearity – this may be the key feature of Deep Learning for images since this paper!
 - use convolution to extract spatial features
 - subsample using spatial average of maps
 - non-linearity in the form of tanh or sigmoids
 - MLP as final classifier
 - sparse connection matrix between layers to avoid large computational cost

LeCun, Y., Bottou, L., Bengio, Y. and Haffner, P., 1998. Gradient-based learning applied to document recognition. Proceedings of the IEEE, 86(11), pp.2278-2324.



LeNet-5 Architecture



- 60K parameters Vs 10–100M parameters (today)
- As you go from left to right n_h and n_w decreases while n_c increases.
- Conv–Pool–Conv–Pool–fc–fc–output (typical architecture)
- Sigmoid/tanh (not ReLU)
- Read sections II and III

PROC. OF THE IEEE, NOVEMBER 1998

1

Gradient-Based Learning Applied to Document Recognition

// LeNet was one of the very first CNNs which helped propel the field of Deep Learning.

Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner

Abstract—

Multilayer Neural Networks trained with the backpropagation algorithm constitute the best example of a successful Gradient-Based Learning technique. Given an appropriate network architecture, Gradient-Based Learning algorithms can be used to synthesize a complex decision surface that can classify high-dimensional patterns such as handwritten characters with minimal error. This paper reviews various methods applied to handwritten character recognition and compares them on a standard handwritten digit recognition task. Convolutional Neural Networks, that are specifically designed to deal with the variability of 2D shapes, are shown to outperform all other techniques.

Real-life document recognition systems are composed of multiple modules including field extraction, segmentation, recognition, and language modeling. A new learning paradigm, called Graph Transformer Networks (GTN), allows such multi-module systems to be trained globally using Gradient-Based methods so as to minimize an overall performance measure.

Two systems for on-line handwriting recognition are described. Experiments demonstrate the advantage of global training, and the flexibility of Graph Transformer Networks.

A Graph Transformer Network for reading bank check is also described. It uses Convolutional Neural Network character recognizers combined with global training techniques to provide record accuracy on business and personal checks. It is deployed commercially and reads several million checks per day.

Keywords— Neural Networks, OCR, Document Recognition, Machine Learning, Gradient-Based Learning, Convolutional Neural Networks, Graph Transformer Networks, Finite State Transducers.

Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. Proceedings of the IEEE, November 1998

57

Other CNNs



1990s to 2012

// Google started in 1998
// Facebook started in 2004

- In the years from late 1990s to early 2010s CNNs were in **incubation**.
- As more and more **data** and **computing power** became available, tasks that CNNs could tackle became more and more interesting.

Birkbeck, University of London 58 © Copyright 2019

58

29

Krizhevsky, A., Sutskever, I. and Hinton, G.E., 2012. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems* (pp. 1097-1105).

ImageNet Classification with Deep Convolutional Neural Networks

Alex Krizhevsky
University of Toronto
kriz@cs.utoronto.ca

Ilya Sutskever
University of Toronto
ilya@cs.utoronto.ca

Geoffrey E. Hinton
University of Toronto
hinton@cs.utoronto.ca

Abstract

 We trained a large, deep convolutional neural network to classify the 1.2 million high-resolution images in the ImageNet LSVRC-2010 contest into the 1000 different classes. On the test data, we achieved top-1 and top-5 error rates of 37.5% and 17.0% which is considerably better than the previous state-of-the-art. The neural network, which has 60 million parameters and 650,000 neurons, consists of five convolutional layers, some of which are followed by max-pooling layers, and three fully-connected layers with a final 1000-way softmax. To make training faster, we used non-saturating neurons and a very efficient GPU implementation of the convolution operation. To reduce overfitting in the fully-connected layers we employed a recently-developed regularization method called “dropout” that proved to be very effective. We also entered a variant of this model in the ILSVRC-2012 competition and achieved a winning top-5 test error rate of 15.3%, compared to 26.2% achieved by the second-best entry.

59

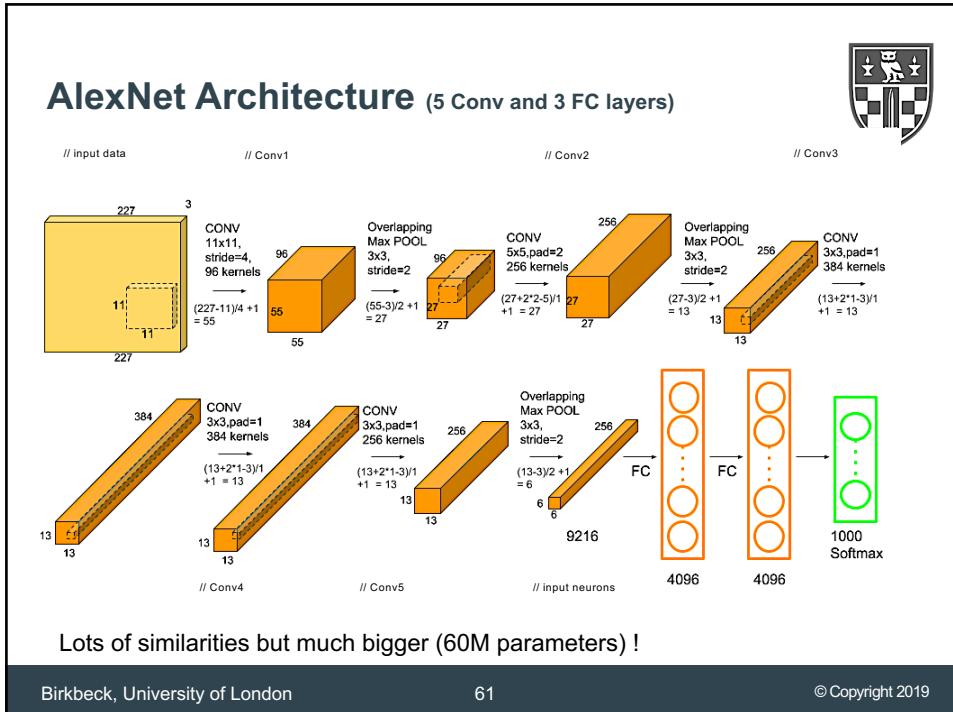
Other CNNs



AlexNet (2012)

- It was a significant breakthrough with respect to the previous approaches and the current widespread application of CNNs can be attributed to this work.
- use of [ReLU](#) as non-linearities
- use of [dropout](#) technique to selectively ignore single neurons during training, a way to avoid overfitting of the model
- [max pooling](#), avoiding the averaging effects of average pooling
- use of multiple [GPUs](#) NVIDIA GTX 580 to reduce training time – 10x faster training time
- Local Response Normalisation – not very useful!
- [Huge impact even beyond computer vision!](#)
- Easier one to read.

60



61

Published as a conference paper at ICLR 2015

VERY DEEP CONVOLUTIONAL NETWORKS FOR LARGE-SCALE IMAGE RECOGNITION

Karen Simonyan* & Andrew Zisserman*
 Visual Geometry Group, Department of Engineering Science, University of Oxford
 {karen, az}@robots.ox.ac.uk

ABSTRACT

In this work we investigate the effect of the convolutional network depth on its accuracy in the large-scale image recognition setting. Our main contribution is a thorough evaluation of networks of increasing depth using an architecture with very small (3×3) convolution filters, which shows that a significant improvement on the prior-art configurations can be achieved by pushing the depth to 16–19 weight layers. These findings were the basis of our ImageNet Challenge 2014 submission, where our team secured the first and the second places in the localisation and classification tracks respectively. We also show that our representations generalise well to other datasets, where they achieve state-of-the-art results. We have made our two best-performing ConvNet models publicly available to facilitate further research on the use of deep visual representations in computer vision.

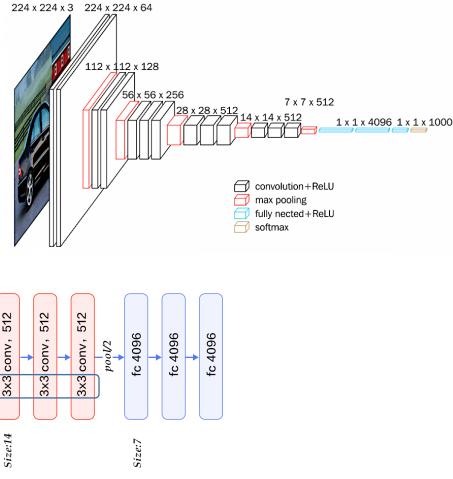
62

Other CNNs



VGG16

- Its main contribution was in showing that the [depth of the network](#) (# of layers) is a critical component for good performance.
- Simplified Architecture
 - $\text{Conv} = 3 \times 3 \text{ filter}$, $s = 1$, same padding
 - $\text{Max-Pool} = 2 \times 2$, $s = 1$
- 138M parameters



Birkbeck, University of London

63

© Copyright 2019

63

Training a single AI model can emit as much carbon as five cars in their lifetimes

Deep learning has a terrible carbon footprint.

by Karen Hao

Jun 6, 2019

The [artificial-intelligence industry](#) is often compared to the oil industry: once mined and refined, data, like oil, can be a highly lucrative commodity. Now it seems the metaphor may extend even further. Like its fossil-fuel counterpart, the process of deep learning has an outsize environmental impact.

In a new paper, researchers at the University of Massachusetts, Amherst, performed a life cycle assessment for training several common large AI models. They found that the process can emit more than 626,000 pounds of carbon dioxide equivalent—nearly five times the lifetime emissions of the average American car (and that includes manufacture of the car itself).

Common carbon footprint benchmarks

in lbs of CO₂ equivalent

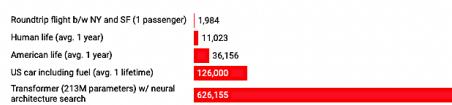


Chart: MIT Technology Review • Source: Strubell et al. • Created with Datawrapper

65

Sabour, S., Frosst, N. and Hinton, G.E., 2017. Dynamic routing between capsules. In *Advances in neural information processing systems* (pp. 3856-3866).

Dynamic Routing Between Capsules

Sara Sabour

Nicholas Frosst

Geoffrey E. Hinton

Google Brain

Toronto

{sasabour, frosst, geoffhinton}@google.com

Abstract

A capsule is a group of neurons whose activity vector represents the instantiation parameters of a specific type of entity such as an object or an object part. We use the length of the activity vector to represent the probability that the entity exists and its orientation to represent the instantiation parameters. Active capsules at one level make predictions, via transformation matrices, for the instantiation parameters of higher-level capsules. When multiple predictions agree, a higher level capsule becomes active. We show that a discriminatively trained, multi-layer capsule system achieves state-of-the-art performance on MNIST and is considerably better than a convolutional net at recognizing highly overlapping digits. To achieve these results we use an iterative routing-by-agreement mechanism: A lower-level capsule prefers to send its output to higher level capsules whose activity vectors have a big scalar product with the prediction coming from the lower-level capsule.

72

Questions?

paul@dcs.bbk.ac.uk

74