Convolutional Neural Networks

Applied Machine Learning, KU Daniel Murnane - May 6th, 2024



8/05/2024

Introduction & Goals

- Nice to meet you!
- I work on ML for high energy physics
 @ NBI
- Goals for today:

UNIVERSITY OF

- Understand how to represent an image
- See how images have traditionally been processed
- Get intuition for the power of a convolution
- Practice "convolution arithmetic"
- Get a feeling for where CNNs fit into the rest of ML
- See a few examples of CNNs in action
- Have borrowed content from <u>MIT Intro to Deep Learning</u>; and Julius Kirkegaard's slides from 2022



Computer Vision



8/05/2024

Convolutional Neural Networks - Daniel Murnane

Modalities in ML

- So far, have seen "tabular data"
- Quick test for "tabularity": Would this work well in Excel?
- The world contains more *modalities* than tabular:
 - Images (today!) i.e. vision
 - Sequences (Wednesday afternoon!) i.e. 1-dimensional time and/or causality
 - Point clouds and graphs (Wednesday morning!) Multi-dimensional time and/or causality
- Most other modalities and senses can be represented as combinations of these:
 - Video = Images + Sequence

UNIVERSITY OF COPENHAGEN

- Robotic action = Tabular + Sequence/Graph
- Audio = Images (waveform) or Sequence (fourier transform)
- Biological & Chemical = Graph (Protein/DNA language models)
- Others might come along... Touch? Smell?
- The hot thing now is combining modalities which is straightforward once you understand how to consume each individually

8/05/2024



Computer Vision

• Today's focus is vision

UNIVERSITY OF COPENHAGEN

- Many diverse tasks: scene reconstruction, object detection, event detection, video tracking, object recognition, 3D pose estimation, learning, indexing, motion estimation, 3D scene modelling, and image restoration
- A beautiful template for the rest of machine learning
- Conventional CV: A huge collection of hand-engineered techniques on a task-by-task basis
- ML CV: A small set of models that can be used for many tasks

8/05/2024



Images as Numbers

- Need way to consume images and apply functions to them
- Most typical way: represent each pixel as a vector
- Greyscale images only need a single value in the vector: the brightness of the pixel
- Colour images need 3 or 4 values in the vector: [R, G, B] or [C, M, Y, K]
- As in **every other ML project**, we normalize those inputs! Conventionally given to us as values [0, 255], therefore divide by 255





A Typical Computer Vision Task

- Let's consider a typical task in CV: classification of the foreground object in an image
- How does a human do this?
- High level feature detection
- We are very sophisticated!
- Imagine if we had to look pixelby-pixel...

High Level Feature Detection

Let's identify key features in each image category



Nose, Eyes, Mouth



Wheels, License Plate, Headlights



Door, Windows, Steps



• Consider a very simple classification problem:

Find images with angle:







- Consider a very simple classification problem:
- Note that we don't care where the feature is *within* the image

Find images with angle:







- Consider a very simple classification problem:
- Note that we don't care where the feature is *within* the image







(color just for visualisation purposes. Think of this a matrix of numbers)

- A "kernel", or a "filter" is a hypothesis about a feature of interest to the task
- We are hand-engineering here, so we can use our intuition
- A kernel is a matrix of the same size of the neighbourhood of pixels we want to "test our hypothesis" on



Define a *kernel, k* =

Place kernel somewhere on image, multiply and sum:



= small number

- A "kernel", or a "filter" is a hypothesis about a feature of interest to the task
- We are hand-engineering here, so we can use our intuition
- A kernel is a matrix of the same size of the neighbourhood of pixels we want to "test our hypothesis" on



Define a *kernel, k* =

Place kernel somewhere on image, multiply and sum:



= small number



- A "kernel", or a "filter" is a hypothesis about a feature of interest to the task
- We are hand-engineering here, so we can use our intuition
- A kernel is a matrix of the same size of the neighbourhood of pixels we want to "test our hypothesis" on



Define a *kernel, k* =

= small number

What happens for a "wrong" image?





- A "kernel", or a "filter" is a hypothesis about a feature of interest to the task
- We are hand-engineering here, so we can use our intuition
- A kernel is a matrix of the same size of the neighbourhood of pixels we want to "test our hypothesis" on



• Applying the filter is actually just a sum over elementwise multiplication:





8/05/2024

• Applying the filter is actually just a sum over elementwise multiplication:





8/05/2024

Max Pooling

- Let's slide the filter across all neighbourhoods in the image
- Keep track of the filter score for each neighbourhood





Max Pooling

- Let's slide the filter across all neighbourhoods in the image
- Keep track of the filter score for each neighbourhood
- Find the maximum score across the whole image





Convolutional Filter

This is precisely what a convolution does!

We have our function! Indeed:

NN(



8/05/2024

= max(conv2D(

- But, we haven't trained anything!
- This is a "hand-crafted" filter:

| 0.1 | 0.1 | 0.9 |
|-----|-----|-----|
| 0.1 | 0.9 | 0.1 |
| 0.9 | 0.1 | 0.1 |

- What if we made the 9 values in the matrix learnable parameters?
- Then we could find the filter that minimises a loss function

Trainable Convolutional Filter







Adding a Latent Space

 Currently everything is "flat" – each pixel multiplied by a single learnable parameter



Image W

• But recall neural networks can learn *many* parameters

Adding a Latent Space

- We can add as much depth to the convolution as we want
- This is a bit hard to understand! There are two levels of matrices...



Adding a Latent Space

- We can add as much depth to the convolution as we want
- This is a bit hard to understand! There are two levels of matrices...



Why not just use a feedforward layer?

- You might ask...
- We have this magical block box: fullyconnected, feedforward neural network
- We know it can approximate any function, provided it is wide and deep enough
- Why not simply feed in the list of pixels as "tabular data" and run through a FFNN?

8/05/2024

• First answer: Translational invariance

| | | | | set. | oduct |
|-------|---|--|----------------------|------------------------|-------------------|
| | 362 | K. Hornik, M. Stinc | hcombe, and H. White | K). | work |
| | | | | acti- | to a |
| | Multilayer Feedforward Nets | | 363 | ther | |
| | | | sΨ | / this | |
| - 1 | | | the | con- | pping |
| | Neural Networks, Vol. 2, pp. 359-366, 1989 Printed in the USA. All rights reserved. | 0693-6680/89 \$3.00 + .00 Copyright © 1989 Pergamon Press plc | s or , | evel. 2.1.) our | aions |
| - 1 | ODICINAL CONTRIBUTION | | | ; no- | |
| - 1 | ORIGINAL CONTRIBUTION | | the | | ching |
| - 1 | | | fol- | | class |
| - 1 | Multilaver Feedfor | ward Networks are | | f and | dfor- |
| - 1 | Universal Ar | nrovimators | | !{x ∈ | layer ars R |
| | Chiversal A | P. Califictoro | alid | (i.e., | out- |
| - 1 | | | DXi- | esults | |
| | Kur H | Iornik | | con- | mean |
| | Technische Ur | siversität Wien | tors | rence | |
| - 1 | | | Π{* | input | |
| - 1 | MAXWELL STINCHCOMB | E AND HALBER WHITE | on- | White | n if it |
| - 1 | University of Cali | fornia, San Diego | oid | (mea- | 1, |
| - 1 | (Received 16 September 1988; rev | rised and accepted 9 March 1989) | | listin- | |
| - 1 | Abstract—This paper rigorously establishes that stando | rd multilayer feedforward networks with as few as one | | ns. tions | Use- |
| | hidden layer using arbitrary squashing functions are ca from one finite dimensional space to another to any o hidden units are available. In this sense, multilayer feedj | pable of approximating any Borel measurable function desired degree of accuracy, provided sufficiently many forward networks are a class of universal approximators. | | next | hres- notes |
| | Keywords—Feedforward networks, Universal approx capability, Stone-Weierstrass Theorem, Squashing func | timation, Mapping networks, Network representation tions, Sigma-Pi networks, Back-propagation networks. | nul- iox- de- | | allant $\pi/2$]) |
| - 1 | 1. INTRODUCTION | any function encountered in applications leads one | ring | ie the | |
| - 1 | It has been nearly twenty years since Minsky and | to wonder about the ultimate capabilities of such networks. Are the successes observed to date re- | ng" dies | | (986) |
| - 1 | Papert (1969) conclusively demonstrated that the simple two-layer perceptron is incapable of usefully | flective of some deep and fundamental approxima- | rise | I only | |
| - 1 | representing or approximating functions outside a | tion capability, or are they merely flukes, resulting from selective reporting and a fortuitous choice of | s of | differ | |
| | very narrow and special class. Although Minsky and Papert left open the possibility that multilayer net- | problems? Are multilayer feedforward networks in | | gale | |
| | works might be capable of better performance, it has | fact inherently limited to approximating only some fairly special class of functions, albeit a class some- | ntal | what | to R |
| | only been in the last several years that researchers have begun to explore the ability of multilaver feed- | what larger than the lowly perceptron? The purpose | esti- | | s |
| | forward networks to approximate general mappings | or this paper is to address these issues. We show that multilayer feedforward networks with as few as one | that | | |
| | trom one finite dimensional space to another. Re- cently, this research has virtually exploded with im- | hidden layer are indeed capable of universal ap- | ren- | n | E N. |
| | pressive successes across a wide variety of applica- | proximation in a very precise and satisfactory sense. Advocates of the virtues of multilayer feedfor- | par- | | |
| | tions. The scope of these applications is too broad to mention useful specifics here: the interested reader | ward networks (e.g., Hecht-Nielsen, 1987) often cite | (sof | → 0. | |
| | is referred to the proceedings of recent IEEE Con- | Kolmogorov's (1957) superposition theorem or its more recent improvements (e.g., Lorentz, 1976) in | nec- | | |
| | ferences on Neural Networks (1987, 1988) for a sam- pling of examples | support of their capabilities. However, these results | Kol- | valent | _ |
| | The apparent ability of sufficiently elaborate feed- | require a <i>different</i> unknown transformation (g in Lorentz's notation) for each continuous function to | ural | In (b) | |
| - 1 | forward networks to approximate quite well nearly | be represented, while specifying an exact upper limit | lden | | |
| | | to the number of intermediate units needed for the | nces | | _ |
| | White's participation was supported by a grant from the Gug- | functions (e.g., logistic, hyperbolic tangent) are used | erty | | |
|) | genheim Foundation and by National Science Foundation Grant SES-8806990. The authors are grateful for helpful suggestions by | in practice, with necessarily little regard for the func- | :h is | | |
| the r | the referees. Requests for reprints should be sent to Halbert White. De- | hidden units increased ad libitum until some desired | | | |
| | and and the state of the state | | | | |



Hornick et al, 1989

Symmetries, Invariances & Inductive Bias

- We *know* that the location of a feature shouldn't matter in an image we can translate a cat in x and y and it's still a cat
- In a CNN, the same convolutional filter is applied in each neighbourhood (aka "shared weights"), so will find a cat in any part of the image: it is translationally invariant
- A FFNN, on the other hand, has a dedicated parameter for *each* pixel it has no concept of "nearness" or a "neighbourhood"
- We have introduced an "inductive bias" by choosing a function (the filter) that is translationally invariant
- A CNN is able to train with *much* less data than a FFNN, and is less prone to overfitting: Because we used our human intuition of translational symmetry



FFNN

Same "Catface" filter



Cat

Different "Catface" neurons



Cat





Convolutional Neural Networks – Daniel Murnane

Convolutional Neural Networks



8/05/2024

Convolutional Neural Networks – Daniel Murnane

Hierarchy and Scale-dependent Features

- Our filter from earlier didn't seem to represent any recognisable object
- But it might be an "edge" or a "corner", which are important low level features in many objects
- Our knowledge of scale is another inductive bias that we can introduce into the CNN architecture





A Typical CNN Architecture

- We can handle scale by having convolutions of increasing size, each one applied to the output of the previous convolution
- Recall the pattern: $MaxPool\left(ReLU\left(Conv2D(x_{ij})\right)\right)$: encoder
- Note: We added a non-linearity here this allows the NN to universally approximate
- We stack multiple encoders to learn features



A Typical CNN Architecture

• We stack multiple encoders to learn features



Convolution Arithmetic

• Let's look at a couple of examples to understand the specifics

CONV2D

O PyTorch

CLASS torch.nn.Conv2d(*in_channels*, *out_channels*, *kernel_size*, *stride=1*, *padding=0*, *dilation=1*, groups=1, bias=True, padding_mode='zeros', device=None, dtype=None) [SOURCE]

 The Conv2D class is simply an implementation of the kernel/filter that we invented earlier, but it has a few bells and whistles that you need to understand



Kernel Size

- Consider the convolution on the right
- Our input image is 4 x 4
- The "kernel size" is the "receptive field" of the convolution: how big is the matrix?
- The kernel size is 3 x 3 (it is conventionally square)





Padding

- Notice in the previous slide that our output is smaller than the input: we don't necessarily want this!
- To give us back our precious pixels, we can add some dummy pixels to the input
- See the example on the right has output size equal to input size





Stride

- Pixels side-by-side are likely to contain very similar information
- It may be a waste of computational resources to consider the neighbourhood of *every single pixel*
- Assuming it's enough that a pixel is in the neighbourhood of *some* convolution, we can set a stride: the number of pixels to slide along

8/05/2024

• Here we stride by 2 pixels

UNIVERSITY OF



Example

- Consider the following:
 - Max Pool of kernel size = 4, stride = 1, padding = 0; applied to
 - Convolution of kernel size = 3, stride = 2, padding = 1; applied to
 - An image of 256 x 256 pixels
- What is the output size?



Back to the Typical Architecture

- Typically, we want to grow our receptive field (by using stride and maxpooling to downsample the image), and deepen the latent space
- This is based on the intuition that there are a small number of complex ideas/features that can fully capture an image



Back to the Typical Architecture

- Okay, so we have convolved and pooled our way to a tiny image with a big latent space
- What happens at the end??
- It depends...



Case Studies



8/05/2024

Convolutional Neural Networks – Daniel Murnane

Classification

- The classic task: classification (or regression)
- As in all classification, we need to transform our $x \times y \times h$ encoded image into a prediction vector (for example, where the entries are class probabilities)
- Simple idea: just line up all the rows of pixels into one big vector (aka "flatten the image") and pass through a FFNN



Classification

- The classic task: classification (or regression)
- As in all classification, we need to transform our $x \times y \times h$ encoded image into a prediction vector (for example, where the entries are class probabilities)
- Simple idea: just line up all the rows of pixels into one big vector (aka "flatten the image") and pass through a FFNN
- But didn't you say this doesn't work very well? Once we have learned features with convolutions and downsampled the image, it works extremely well. Doesn't work well on the raw input pixels



Classification: Breast Cancer Screening



CNN-based system outperformed expert radiologists at detecting breast cancer from mammograms



Breast cancer case missed by radiologist but detected by Al



Instance and Semantic Segmentation

UNIVERSITY OF COPENHAGEN

8/05/2024

 What if we want to label all the objects in an image, and find their location? 128 64 64 2 • This is segmentation: rather than a classification for the image, give a classification for *each pixel* input output image 🔶 segmentation tile map • But this requires returning an output the same size as the input image 128 128 • Enter the "Unet": allows a large receptive field, while still returning the original dimensions Contains a "transposed convolution" 256 conv 3x3, ReLU copy and crop 512 max pool 2x2 up-conv 2x2 conv 1x1

Semantic Segmentation: Tritium Detection

- Given the fallout of a nuclear disaster, can you detect radioactive decays from soil in a silicon charge-couple device?
- We hope so! This is a semantic segmentation problem:



Generative Models

- Recall the Unet: it goes all the way down to a single "pixel" of very large latent space
- Called the "bottleneck" used to capture information about the whole image
- The grey arrows are "residual connections": simply concatenate/sum two latent spaces together! Simple but powerful

8/05/2024

UNIVERSITY OF COPENHAGEN





Generative Models

- Recall the Unet: it goes all the way down to a single "pixel" of very large latent space
- Called the "bottleneck" used to capture information about the whole image
- We can train a network that *starts* at the bottleneck to perform transposed convolutions and *generate* an image
- Many different ways to train, but they almost always rely on bottleneck→transposed convolution





Generative Models

 CycleGAN is a sophisticated sequence of transposed convolutions and bottlenecks





Transfer Learning with Feature Extraction

 What if we spend \$1million training a cat/dog classifier, then realise we actually want to classify lions and wolves – or worse: aeroplanes, boats and cars





Transfer Learning with Feature Extraction

- What if we spend \$1million training a cat/dog classifier, then realise we actually want to classify lions and wolves or worse: aeroplanes, boats and cars
- Intuition: the features learned from a cat picture (edges for example) are relevant for other image classifiers
- We can slice off the final task-specific layer of a model and still use its feature layers



Transfer Learning: ImageNet

- ImageNet has 14 million images, with 20,000 (!!!) classes
- There are many publicly available CNNs trained on ImageNet that have feature layers for transfer learning
- ResNet is a classic model that combines convolutions and residual connections

Image Classification on ImageNet



import torch
model = torch.hub.load('pytorch/vision:v0.10.0', 'resnet18', pretrained=True)



8/05/2024

Other Dimensions of Convolution: Sequences & Videos

- We have only talked about 2D convolutions, but \rightarrow
- If we have a sequence of scalar measurements (e.g. histogram), where nearby measurements should be related somehow, then we can use a 1D convolution
- If we have a sequence of vector measurements, then we can use a 2D convolution
- If we have a sequence of 2D measurements (e.g. video), then we can use a 3D convolution
- In principle, could go ever higher!

8/05/2024



| | O PyTorch |
|---|--------------------|
| | Convolution Layers |
| Ì | nn.Conv1d |
| n | nn.Conv2d |
| | nn.Conv3d |

Other Dimensions of Convolution: Sequences & Videos

- We can perform a 3D convolution across [*t*, *x*, *y*]
- That is, instead of:



Image W

- We need to use a $t \times x \times y \times h$ tensor to multiply each time-pixel
- This might capture concepts that are not static, e.g. emotions...





Data Augmentation & Learned Symmetries

- We know we have translational invariance (the sliding kernel!)
- Do we have rotational invariance? Dilation invariance? <u>Lorentz invariance</u>?
- The MIT course certainly seems to think so...



Image is represented as matrix of pixel values... and computers are literal! We want to be able to classify an X as an X even if it's shifted, shrunk, rotated, deformed.

8/05/2024

UNIVERSITY OF COPENHAGEN



Data Augmentation & Learned Symmetries

- In fact, we don't have these other symmetries "hard-coded" in, like we do with translation
- If we want to include them, we need to learn them: Enter data augmentation
- Principle: In training, randomly transform image samples with rotations, noise, dilations, colour changes, etc.
- Practice: Use a Transform function

8/05/2024







https://arxiv.org/pdf/2012.09699.pdf

Key Take-aways

- The convolution is a building block in <u>almost any architecture</u> applied to vision and video
- Convolution arithmetic is kind of a pain, but that's what print(tensor.shape) statements are for
- Convolutions can also be applied to sequential data, can learn symmetries in augmented data, and can segment images
- These days: Always start with a pre-trained model it's a free lunch!
- The humble CNN proves the value of many ideas in ML: symmetry, hierarchy, inductive bias, data augmentation, generative models, etc.



CNN FAQ



8/05/2024

Convolutional Neural Networks – Daniel Murnane

• Recall our toy model





• Recall our toy model

COPENHAGEN



The window ("kernel size") is 3x3, but the input shape (in each pixel) is just one number, and the output shape (for each pixel) is just one number

Filter

• Recall our toy model

Image neighbourhood

COPENHAGEN



The window ("kernel size") is 3x3, but the input shape (in each pixel) is just one number, and the output shape (in each pixel) is just one number

- To make a CNN more powerful, we need to be able to take in any number of features (e.g. color) and we want *many* filter channels in a convolution. That is, we need input length and output length
- Our general convolution is thus

$$\sum_{ij} x_{ij}^M W_{ij}^{MN} = \sum_{ij} m_{ij}^N = h^N$$

where x_{ij}^{M} is the input window, W_{ij}^{MN} is the learnable kernel with the ij^{th} entry an $M \times N$ matrix. M is input shape (e.g. 3 colors), N is output shape, i is the window rows, j is the window columns



When do I do max pooling?

- Max pooling and the convolutional kernel are **completely separate**, and we don't even need to do pooling
- They both use a "window" system, of looking at neighborhoods of pixels, but the window size of the kernel does not have to be the same size as the max pooling window size
- Max pooling is just one choice of pooling. Any kind of aggregation will work: mean, sum, min, or even more exotic kinds like variance



But does the pooling reduce the image size?

- If we pad our image on the border, and slide the kernel or pooling window one pixel at a time, our output array will be the same shape as the input array: the kernel and pooling do not inherently downsample an image
- It is only when we increase the stride (the number of pixels we move the window each time) that the output shape gets reduced
- Either kernel and/or pooling use windows, so they can both use stride, and thus can both be used to down-sample an image. Which one you use to down-sample is part intuition, part trial-and-error

