Topics:

Convolutional Neural Networks

CS 4644-DL / 7643-A ZSOLT KIRA

- Assignment 1 Due Thursday!!!
 - DO NOT SEARCH FOR CODE!!!!
- Assignment 2
 - Implement convolutional neural networks
 - From scratch
 - Using Pytorch
- GPU resources
 - For assignments, can use CPU or Google Colab
 - Projects:
 - Google Cloud Credits
 - PACE-ICE

W5: June 11	Project planning session Project Proposal Due June 15th 11:59pm (grace period until June 17th) CNNs notes, CNNs Backprop notes.
W6: June 16	CNNs continued, Regularization, Augmentation, Transfer Learning READING: A ConvNet for the 2020s

Please do the reading!

The connectivity in linear layers doesn't always make sense



How many parameters?

M*N (weights) + N (bias)

Hundreds of millions of parameters **for just one layer**

More parameters => More data needed

Is this necessary?



Limitation of Linear Layers



Image features are spatially localized!

- Smaller features repeated across the image
 - Edges
 - Color
 - Motifs (corners, etc.)
- No reason to believe one feature tends to appear in one location vs. another (stationarity)

Locality of Features

Can we induce a *bias* in the design of a neural network layer to reflect this?









We can learn **many** such features for this one layer

- Weights are **not** shared across different feature extractors
- Parameters: (K₁×K₂ + 1) * M where M is number of features we want to learn



Idea 3: Learn Many Features



This operation is extremely common in electrical/computer engineering!

w(t) y(t)XAY $(z) = e^{-\left(\frac{z}{2}-\frac{z}{2}\right)}$ $Y(z) = (X_{a} + w)(z)$ = $a = -\infty \times (z - a) w(a) da$ = $(4 + x)(z) = 3 \times (a) w(z) da$



From https://en.wikipedia.org/wiki/Convolution





We will make this convolution operation a layer in the neural network

- Initialize kernel values randomly and optimize them!
- These are our parameters (plus a bias term per filter)









 $y(0,0) = x(-2,-2)k(2,2) + x(-2,-1)k(2,1) + x(-2,0)k(2,0) + x(-2,1)k(2,-1) + x(-2,2)k(2,-2) + \dots$

Mathematics of Discrete 2D Convolution



1. Flip kernel (rotate 180 degrees)



2. Stride along image





== Cross-Correlation with Flipped Kernel



$$y(r,c) = (x * k)(r,c) = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} x(r+a,c+b) k(a,b)$$

(0,0)



Since we will be learning these kernels, this change does not matter!





$$X(0:2,0:2) = \begin{bmatrix} 200 & 150 & 150 \\ 100 & 50 & 100 \\ 25 & 25 & 10 \end{bmatrix} \qquad K' = \begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix} \longrightarrow X(0:2,0:2) \cdot K' = 65 + \text{bias}$$

Dot product (element-wise multiply and sum)





































Convolution and Cross-Correlation











Why Bother with Convolutions?

Convolutions are just **simple linear operations**

Why bother with this and not just say it's a linear layer with small receptive field?

- There is a **duality** between them during backpropagation
- Convolutions have various mathematical properties people care about

This is historically how it was inspired





Input & Output Sizes



Convolution Layer Hyper-Parameters

Parameters

- in_channels (int) Number of channels in the input image
- out_channels (int) Number of channels produced by the convolution
- kernel_size (int or tuple) Size of the convolving kernel
- stride (int or tuple, optional) Stride of the convolution. Default: 1
- padding (int or tuple, optional) Zero-padding added to both sides of the input. Default: 0
- padding_mode (string, optional) 'zeros', 'reflect', 'replicate' or 'circular'. Default: 'zeros'

Convolution operations have several hyper-parameters

From: https://pytorch.org/docs/stable/generated/torch.nn.Conv2d.html#torch.nn.Conv

Output size of vanilla convolution operation is $(H - k_1 + 1) \times (W - k_2 + 1)$

This is called a "valid" convolution and only applies kernel within image





We can **pad the images** to make the output the same size:

Zeros, mirrored image, etc.

Note padding often refers to pixels added to one size (P = 1 here)







 $W + 2 - k_2 + 1$





We can move the filter along the image using larger steps (stride)

- This can potentially result in loss of information
- Can be used for dimensionality reduction (not recommended)

Stride = 2 (every other pixel)









Stride can result in **skipped pixels**, e.g. stride of 3 for 5x5 input



W





We have shown inputs as a **one-channel image** but in reality they have three channels (red, green, blue)

In such cases, we have 3-channel kernels!





We have shown inputs as a **one-channel image** but in reality they have three channels (red, green, blue)

In such cases, we have 3-channel kernels!



Similar to before, we perform **element-wise multiplication** between kernel and image patch, summing them up (dot product)

Except with $k_1 * k_2 * 3$ values





We can have multiple kernels per layer

We stack the feature maps together at the output

Number of channels in output is equal to *number* of kernels





Number of parameters with N filters is: $N * (k_1 * k_2 * 3 + 1)$

Number of Parameters





Just as before, in practice we can vectorize this operation

Step 1: Lay out image patches in vector form (note can overlap!)

Input Image



Adapted from: https://petewarden.com/2015/04/20/why-gemm-is-at-the-heart-of-deep-learnin_/





Just as before, in practice we can vectorize this operation

Step 2: Multiple patches by kernels

Input Matrix

Kernel Matrix



Adapted from: https://petewarden.com/2015/04/20/why-gemm-is-at-the-heart-of-deep-learning/





Backwards Pass for Convolution Layer



It is instructive to calculate **the backwards pass** of a convolution layer

- Similar to fully connected layer, will be simple vectorized linear algebra operation!
- We will see a **duality** between cross-correlation and convolution





Backwards Pass for Conv Layers

$$y(r,c) = (x * k)(r,c) = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} x(r+a,c+b) k(a,b)$$

(0,0)



 $W = 5 \quad (H-1, W-1)$





$$y(r,c) = (x * k)(r,c) = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} x(r+a,c+b) k(a,b)$$



Some simplification: 1 channel input, 1 kernel (channel output), padding (here 2 pixels on right/bottom) to make output the same size





$$y(r,c) = (x * k)(r,c) = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} x(r+a,c+b) k(a,b)$$

 $|y| = H \times W$

$\frac{\partial L}{\partial y}$? Assume size $H \times W$ (add padding)

$$\frac{\partial L}{\partial y(r,c)}$$
 to access element







Backpropagation Chain Rule

Georgia Tech∦

Gradient for Convolution Layer


$\frac{\partial L}{\partial k} = \frac{\partial L}{\partial h^{\ell}} \quad \frac{\partial h^{\ell}}{\partial k}$

Gradient for weight update

Calculate one pixel at a time

$$\frac{\partial L}{\partial k(a,b)} = \frac{\partial L}{\partial h^{\ell}} \quad \frac{\partial h^{\ell}}{\partial k(a,b)}$$

What does this weight affect at the output?

Everything!



 $W = 5 \qquad (H-1, W-1)$

What a Kernel Pixel Affects at Output



Need to incorporate all upstream gradients:

 $\left\{\frac{\partial L}{\partial y(0,0)}, \frac{\partial L}{\partial y(0,1)}, \dots, \frac{\partial L}{\partial y(H,W)}\right\}$

Chain Rule: $\frac{\partial L}{\partial k(a,b)} = \sum_{r=0}^{H-1} \sum_{c=0}^{W-1} \frac{\partial L}{\partial y(r,c)} \frac{\partial y(r,c)}{\partial k(a,b)}$ Sum over Upstream We will all output gradient compute pixels (known)





 $W=5 \qquad (H-1,W-1)$





 $\frac{\partial y(r,c)}{\partial k(a,b)} =?$





For output at y(r,c), where is placement of kernel?







Chain Rule over all Output Pixels



 $\frac{\partial y(r,c)}{\partial k(a,b)} =?$



Η

Which x pixel is multiplied by k(a,b)



W





W



 $\frac{\partial y(r,c)}{\partial k(a,b)}$ =?

Reasoning:

- Cross-correlation is just "dot product" of kernel and input patch (weighted sum)
- When at pixel y(r, c), kernel is on input x such that k(0, 0) is multiplied by x(r, c)
- But we want derivative w.r.t. k(a, b)
 - k(0,0) * x(r,c), k(1,1) * x(r+1,c+1), k(2,2) * x(r+2,c+2)
 - => in general k(a, b) * x(r + a, c + b)
 - Just like before in fully connected layer, partial derivative w.r.t. k(a, b) only has this term (other x terms go away because not multiplied by k(a, b)).









 $\frac{\partial y(r,c)}{\partial k(a,b)} = x(r+a,c+b)$

$$\frac{\partial L}{\partial k(a,b)} = \sum_{r=0}^{H-1} \sum_{c=0}^{W-1} \frac{\partial L}{\partial y(r,c)} x(r+a,c+b)$$

Does this look familiar?

Cross-correlation between upstream gradient and input! (until $k_1 \times k_2$ output)







Forward Pass







 $\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \quad \frac{\partial y}{\partial x}$

Gradient for input (to pass to prior layer)

Calculate one pixel at a time

$$\frac{\partial L}{\partial x(r',c')}$$

What does this input pixel affect at the output?

Neighborhood around it (where part of the kernel touches it)

(0,0)



 $W = 5 \qquad (H-1, W-1)$



What an Input Pixel Affects at Output



W





W













This is where the corresponding locations are for the **output**









Summing Gradient Contributions







Summing Gradient Contributions







Let's derive it analytically this time (as opposed to visually)





Summing Gradient Contributions



Definition of cross-correlation (use a', b' to distinguish from prior variables):

$$y(r',c') = (x * k)(r',c') = \sum_{a'=0}^{k_1-1} \sum_{b'=0}^{k_2-1} x(r' + a',c' + b') k(a',b')$$

Plug in what we actually wanted :

$$y(r'-a,c'-b) = (x * k)(r',c') = \sum_{a'=0}^{k_1-1} \sum_{b'=0}^{k_2-1} x(r'-a+a',c'-b+b') k(a',b')$$

What is
$$\frac{\partial y(r'-a,c'-b)}{\partial x(r',c')} = k(a,b)$$
 (we want term with $x(r',c')$ in it;
this happens when $a' = a$ and $b' = b$)



Calculating the Gradient

Plugging in to earlier equation:

$$\frac{\partial L}{\partial x(r',c')} = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} \frac{\partial L}{\partial y(r'-a,c'-b)} \frac{\partial y(r'-a,c'-b)}{\partial x(r',c')}$$

Does this look familiar?

$$=\sum_{a=0}^{k_1-1}\sum_{b=0}^{k_2-1}\frac{\partial L}{\partial y(r'-a,c'-b)}k(a,b)$$

Again, all operations can be implemented via matrix multiplications (same as FC layer)! Convolution between upstream gradient and kernel!

(can implement by flipping kernel and cross- correlation)





- Convolutions are mathematical descriptions of striding linear operation
- In practice, we implement **cross-correlation neural networks!** (still called convolutional neural networks due to history)
 - Can connect to convolutions via duality (flipping kernel)
 - Convolution formulation has mathematical properties explored in ECE
- Duality for forwards and backwards:
 - Forward: Cross-correlation
 - Backwards w.r.t. K: Cross-correlation b/w upstream gradient and input
 - Backwards w.r.t. X: Convolution b/w upstream gradient and kernel
 - In practice implement via cross-correlation and flipped kernel
- All operations still implemented via **efficient linear algebra** (e.g. matrixmatrix multiplication)





Pooling Layers



- Dimensionality reduction is an important aspect of machine learning
- Can we make a layer to explicitly down-sample image or feature maps?



Yes! We call one class of these operations pooling operations

Parameters

- kernel_size the size of the window to take a max over
- stride the stride of the window. Default value is kernel_size
- padding implicit zero padding to be added on both sides

From: https://pytorch.org/docs/stable/generated/torch.nn.MaxPool2d.html#torch.nn.MaxPool2/





Example: Max pooling

Stride window across image but perform per-patch max operation

 $X(0:2,0:2) = \begin{bmatrix} 200 & 150 & 150 \\ 100 & 50 & 100 \\ 25 & 25 & 10 \end{bmatrix} \implies \max(0:2,0:2) = 200$



Max Pooling

How many learned parameters does this layer have?





- makes the representations spatially smaller
- saves computation (GPU mem & speed), allows go deeper
- operates over each activation map independently:



From: Slides by CS 231n, Danfei Xu



Since the **output** of convolution and pooling layers are **(multi-channel) images**, we can sequence them just as any other layer







This combination adds some invariance to translation of the features

If feature (such as beak) translated a little bit, output values still remain the same







Convolution by itself has the property of equivariance

If feature (such as beak) translated a little bit, output values move by the same translation





W = 5





Simple Convolutional Neural Networks



Since the **output** of convolution and pooling layers are **(multi-channel) images**, we can sequence them just as any other layer









Alternating Convolution and Pooling





Geo







Receptive Fields



These architectures have existed **since 1980s**



Image Credit: Yann LeCun, Kevin Murphy



Georg a Tech

Handwriting Recognition



Image Credit: Yann LeCun Georg a

Translation Equivariance (Conv Layers) & Invariance (Output)



Image Credit: Yann LeCun Georgia

(Some) Rotation Invariance



Image Credit: Yann LeCun Georga

(Some) Scale Invariance



Image Credit: Yann LeCun Georgaa





