

CNNs

John M



- ① Get data ready
 - ② Build or pick up a pretrained model
to suit your problem
 - ③ Fit the model to the data and make a prediction
 - ④ Evaluate the model
 - ⑤ Improve through experimentation
 - ⑥ Save and reload your trained model
- Build a training loop ↗
pick a loss function & optimizer ↘

SHIFT + CMD + SPACE to read the docstrings

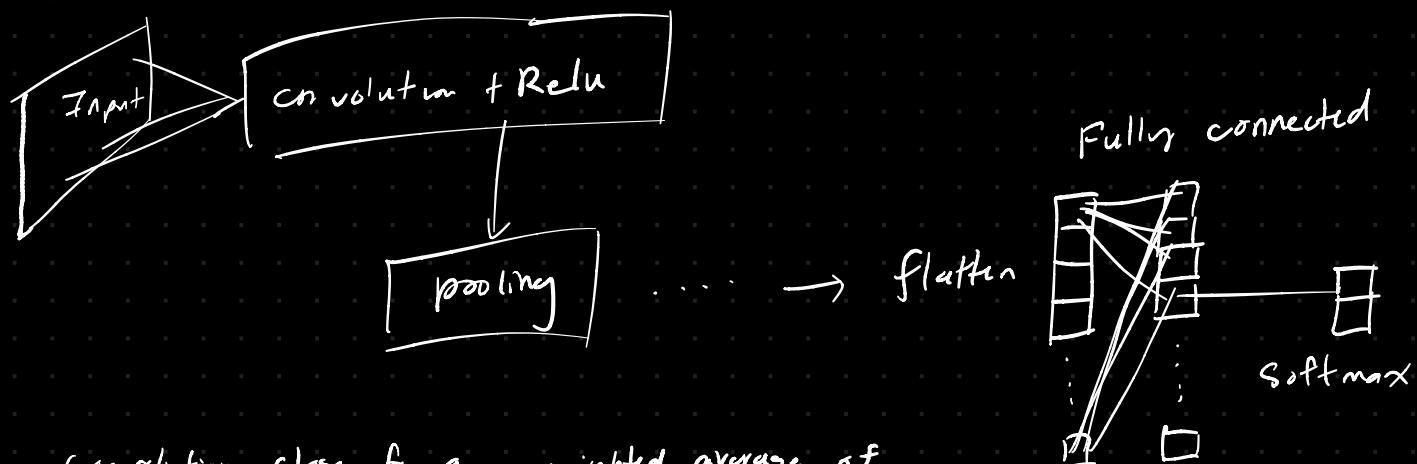
CNN Slideshow -03/17/24

A CNN is a type of neural network that assumes the input is an image.

The layers of a CNN have neurons arranged in 3 dimensions

- ① height
- ② width
- ③ depth

Convolutions use a kernel (matrix with learnable weights, the kernel is then convolved with neurons in that layer)



A convolution takes a weighted average of values in a specific region.

$$g(x, y) = w \cdot f(x, y) = \sum_{i=-a}^a \sum_{j=-b}^b w(i, j) f(x-i, y-j)$$

w is the filter kernel

All elements are considered by $i \in [-a, a]$
 $j \in [-b, b]$

Hyper parameters

Depth - # of filters

Stride - # of pixels to slide the kernel

Zero-padding → controls output size

pooling layers

Max pooling → outputs the maximum in the area

pooling reduces # of parameters and the training cost

Non-linearity layers

Sigmoid

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

$\tanh(x)$

$\tanh(x)$

ReLU

$$\max(0, x)$$

Leaky ReLU

$$\max(0.1x, x)$$

Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\alpha(x) = \begin{cases} x & x \geq 0 \\ e^x - 1 & x < 0 \end{cases}$$

Fully connected layer

The last layer of a CNN

Full connections to all activations in the previous layer.

Softmax → normalizes output to be between 0 and 1.

→ represents the probability that the input belongs to a specific class.

Input \rightarrow neuron \rightarrow dot product \rightarrow non-linearity

out put \rightarrow single differentiable score / classes score

Loss functions in a CNN \rightarrow SVM + Soft max
(Support vector machine)

There is the explicit assumption that the input is an image \rightarrow reduce the # of parameters in the network.

Regular Neural Network

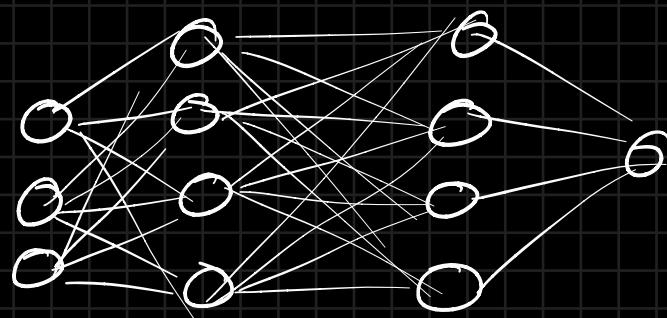
Neural networks receive a single input

and then this input is passed onto the "hidden" layers of the network \rightarrow transforming the output in the process.

Each layer consists of neurons, that function independently.

Each neuron in each layer is connected fully to the neurons of the previous layer.

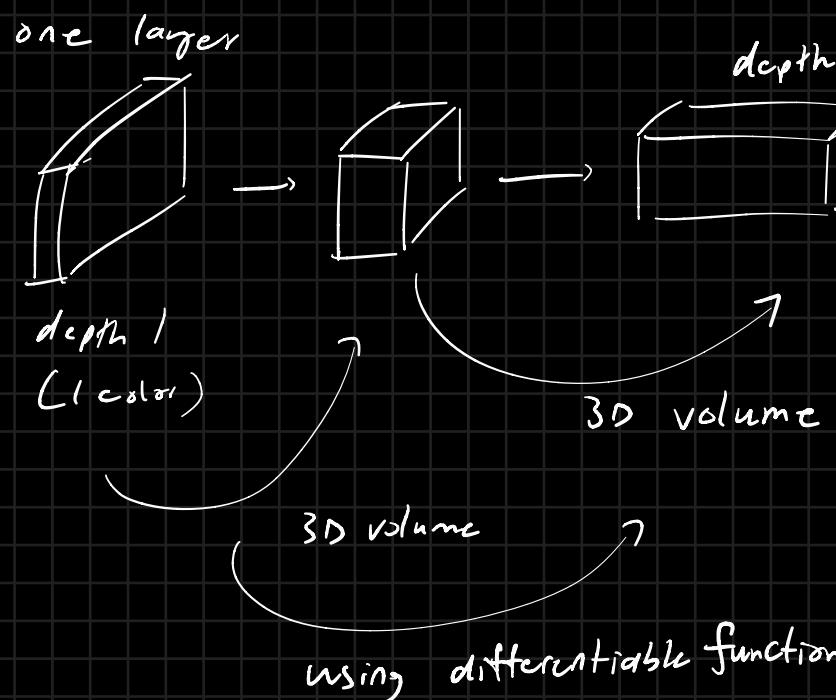
The last layer is called the output layer and allows the user to see class scores



Regular neural networks require too many weights if the input size is a large image.

→ Also don't want a huge # of parameters to avoid overfitting

CNN architecture is constrained as input has a height, weight, and depth
Instead of full connections between hidden layers, neurons are connected to only a small subset of the previous layer's neurons.



3 types of layers in a CNN:

convolutional layers → pooling layer → fully connected layer

Input → conv → pool → FC

Input holds the original pixel data

conv layer will take in a subset of input and output corresponding values

ReLU applies $\max(0, x)$ to each neuron

pooling performs a down sampling operation along the spatial dimensions along spatial dimensions, Scaling down

FC layer computes the class scores and has height 1, width 1, and depth # of classes.

conv / FC perform transformations of values using weights and biases

parameters with ^{conv} / FC are updated using gradient descent.

Input $32 \times 32 \times 3$

conv $32 \times 32 \times 12$

ReLU $32 \times 32 \times 12$

pool $16 \times 16 \times 12$

FC $1 \times 1 \times 10$

Conv / FC → function of weights & biases

ReLU / pool → fixed function

input volume → output volume

Conv layer: Set of learnable filters

Filters \Rightarrow extend RGB height, but less width and length

e.g. $5 \times 5 \times 3$

Slide / convolve each filter across the width and height of the input volume and compute dot products between overlapping entries at any position of the filter.

\Rightarrow produces 2 dimension activation map / matrix that gives responses of that filter at every position

The goal is that through training, the filters learn to activate when they see some visual feature.

With a set of filters, each produce a separate 2 dimensional activation map / matrix

Stacking along dimensions of height, to produce the output 3D matrix

as a brain analogy, every output in the output layer is a neuron that only looks at a small region of the image and shares learnable parameters with all neurons to the left and right spatially

Small subset of image \rightarrow hyperparameter
aka receptive field / filter size

Connectivity along the depth = depth of input volume (RGB)

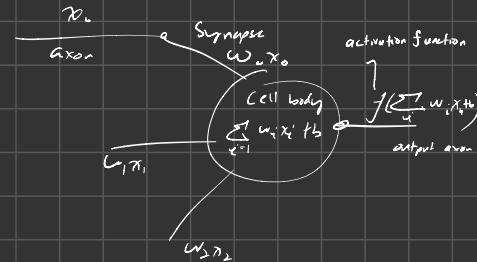
Each filter has one bias only *

Three hyperparameters control the size of the output volume : Depth, stride, and zero padding

Depth of the output volume (# of filters, each learning to look for something different in the input)
Neurons looking at the same region are called a depth column

Stride : How far we slide the filter.

padding : zeros around the border to control spatial size of the output



Not always RGB
simply height
of input
volume

Spatial size of output volume

Input Volume (w)

Stride must be picked to make

Field / filter size (F)

Stride (s)

sure the convolution is symmetric

Zero padding (p)

Since the convolution / filter is a

$$\frac{w - F + 2p}{s} + 1$$

Square wrt its height and width,

.....
0 .. 0 .. 0 .. 0 .. 0 ..
0 .. 0 .. 0 .. 0 .. 0 ..
0 .. 0 .. 0 .. 0 .. 0 ..
0 .. 0 .. 0 .. 0 .. 0 ..
0 .. 0 .. 0 .. 0 .. 0 ..
0 .. 0 .. 0 .. 0 .. 0 ..

$$\frac{w - F + 2p}{s} + 1 \text{ must be divisible by } F$$

If we set zero padding to be $p = \frac{F-1}{2}$ when

Stride = 1, then, the input & output volume have the same sizes

chosen stride must allow for $\frac{w - F + 2p}{s} + 1$ to be an integer

If one feature is useful to compute at some spatial position x, y , then it

Should also be useful to compute at different position x_2, y_2

denote a depth slice (slice of height)

constraint neurons in each depth slice to use the same
weights and bias. Set of weights, one bias \Rightarrow per slice/filter

During back propagation, every neuron in the volume will compute the gradient for its weights but these gradients will be added up across each depth slice and only update a single set of weights per slice.

If all neurons in a depth slice use the same weights vector, then the forward pass of the convolutional layer is computed as a convolution of neurons weights with the input volume.

In Numpy,

Input x num array

A depth column at x_{iy} would be the activations $x(x, y, :)$

A depth slice of depth d would be the activations $x[:, :, d]$

Conv layer example

Input x has shape $(11, 11, 4)$

If we use 0 padding and a filter size of 5
with stride 2

$$\text{Output volume: } \frac{11-5}{2+1} = 4$$

Activation map:

$$\begin{aligned} v[0, 0, 0] &= np.sum(x[0:5, 0:5, :] * w_0) + b_0 \\ v[1, 0, 0] &= np.sum(x[2:7, 0:5, :] * w_0) + b_0 \\ v[2, 0, 0] &= np.sum(x[4:9, 0:5, :] * w_0) + b_0 \\ v[3, 0, 0] &= np.sum(x[6:11, 0:5, :] * w_0) + b_0 \end{aligned} \quad \left. \begin{array}{l} \text{element wise} \\ \text{Same } w_0 \text{ and} \\ \text{so (parameter sharing)} \end{array} \right\}$$

$$w_0.shape = (5, 5, 4)$$

Computing dot product

Conv layer

Accepts a volume of size $W_1 \times H_1 \times D_1$

Requires four hyper parameters

The filters F

Their spatial extent F

The stride S

The amount of zero padding P

output volume $W_2 \times H_2 \times D_2$ where

$$W_2 = ((W_1 - F + 2P) / S) + 1$$

$$H_2 = ((H_1 - F + 2P) / S) + 1$$

$$D_2 = k$$

With parameter sharing, $F \times F \times D$, weights per filter for a total of

$F \times F \times D \times k$ weights, and k biases

The d th depth slice is the result of performing a valid convolution
with stride S , offset by d th bias

Common setting: $F=3$, $S=1$, $P=1$

weights are shared across a slice, per filter

①

②

if and only if

otherwise, weights are different

Biases are shared across a slice, per filter.

As matrix multiplication:

local regions are stretched out into columns in an operation called im2col
=> column vector

Ex. input = $227 \times 227 \times 3$ $11 \times 11 \times 3 \approx 363$

Kernel = $11 \times 11 \times 3$ at stride 4

$$\frac{227-11}{4+1} = 55 \text{ locations along both width and length}$$

output matrix => (363×3025)
 $\times \text{-col}$

$(96 \times 55 \times 55)$

Stretch out weights: 96 filters eg: (96×363)

↑ reshape

one convolution: perform np.dot (w-row, x-col) => (96×3025)

Back propagation : Spatially flipped filters

Pooling layer \Rightarrow reduces spatial size

Resized using the max operation

Common: size 2×2 filter, with stride 2 at every depth slice.

For input $W_1 \times H_1 \times D_1$,

and spatial extent F

Stride S

produced a volume of $CW_2 \times H_2 \times D_2$

Sometimes pooling is not used
and larger stride replaces
pooling to prevent over fitting

$$W_2 = (W_1 - F) / S + 1$$

$$H_2 = (H_1 - F) / S + 1$$

$$D_2 = D_1$$

Can also perform average pooling

Back propagation \Rightarrow routing gradients that had the highest value in the forward pass. others are just zero.

Fully connected

full connections to all activations in the previous layer as seen in regular neural networks.

\Rightarrow Matrix multiplication with bias offset

For any conv layer, an FC matrix exists that can implement the same forward function (Most entries would be zero)

Conversely, if we had a FC layer looking at $7 \times 7 \times 512$ input, and $k=4096$, the conv layer would be $F=7$, $p=0$, $S=1$, $K=4096$

Filter Size = input volume

output: $1 \times 1 \times 4096$

only 1 depth fits

e.g. $224 \times 224 \times 3 \Rightarrow 7 \times 7 \times 512$

5 pooling layers

$F=1 \rightarrow 1 \times 1 \times 1000$
 $k=100$

FC layers work and loss

FC \rightarrow conv: $7 \times 7 \times 512 \Rightarrow F=7$ (above) $\Rightarrow 1 \times 1 \times 4096$
 $k=4096$

If we have a larger image and run
a convolution,

our class output is not 1D any more

\Rightarrow Take the average across extra dimensions at all

Spatial positions (as long as the initial convolution is legal)

Input \rightarrow $(\text{Conv} \rightarrow \text{ReLU})^n \rightarrow \text{pool} \times m \rightarrow (\text{FC} \rightarrow \text{ReLU})^k \rightarrow \text{FC}$

$$n \geq 0 \quad m \geq 0 \quad k \geq 0 \quad k < 3$$

Input $\rightarrow (\text{Conv} \rightarrow \text{ReLU} \rightarrow \text{pool}) \times 2 \rightarrow \text{FC} \rightarrow \text{ReLU} \rightarrow \text{FC}$

Stacks of conv filters are better than one large conv layer

More expressive with fewer parameters

CNNs from Scratch

Convolve vs Correlate

$$I * K = I \star \text{rot 180}(K)$$

3 dim block of data \rightarrow depth = 3

Kernel 1 (3 dim overall)

$$\begin{bmatrix} K_{11} \\ K_{12} \\ K_{13} \end{bmatrix} \quad \begin{bmatrix} B_1 \end{bmatrix}$$

Kernel 2 (3 dim)

$$\begin{bmatrix} K_{21} \\ K_{22} \\ K_{23} \end{bmatrix} \quad \begin{bmatrix} B_2 \end{bmatrix}$$

output

$$\begin{bmatrix} Y_1 \end{bmatrix}$$

biases

$$\begin{bmatrix} Y_2 \end{bmatrix}$$

Cross correlation

$$Y_1 = \beta_1 + \gamma_1 * k_{11} + \gamma_2 * k_{12} + \gamma_3 * k_{13} \dots x_n k_{1n}$$

$$Y_2 = \beta_2 + \gamma_1 * k_{21} + \gamma_2 * k_{22} + \gamma_3 * k_{23}$$

:

$$Y_d = \beta + \gamma_1 * k_{d1} + \gamma_2 * k_{d2} + \gamma_3 * k_{d3} \dots x_n k_{dn}$$

Forward propagation

$$Y_i = \beta_i + \sum_{j=1}^n x_j * k_{ij}, i = 1 \dots d$$

$$Y = \underbrace{\beta + K \cdot \underbrace{x}_{\text{cross correlation}}}_{\text{dense layer}}$$

Back propagation

$$\frac{\partial C}{\partial Y_i} = \frac{\partial C}{\partial K_{ij}} \quad \text{①}$$

$$\frac{\partial C}{\partial B_i} \quad \text{②}$$

$$\frac{\partial C}{\partial X_j} \quad \text{③}$$

$$Y_i = B_i + \sum_{j=1}^n X_j * K_{ij}$$

Take one layer, and expand all convolutions

$$\frac{\partial E}{\partial K_{ii}} = \frac{\partial C}{\partial Y_{11}} \frac{\partial Y_{11}}{\partial K_{11}} + \frac{\partial C}{\partial Y_{12}} \frac{\partial Y_{12}}{\partial K_{11}} + \frac{\partial C}{\partial Y_{21}} \frac{\partial Y_{21}}{\partial K_{11}} + \frac{\partial C}{\partial Y_{22}} \frac{\partial Y_{22}}{\partial K_{11}}$$

$$\begin{aligned} \frac{\partial C}{\partial K_{11}} &= \frac{\partial C}{\partial Y_{11}} \\ \frac{\partial C}{\partial K_{12}} &= \frac{\partial C}{\partial Y_{12}} \\ \frac{\partial C}{\partial K_{21}} &= \frac{\partial C}{\partial Y_{21}} \\ \frac{\partial C}{\partial K_{22}} &= \frac{\partial C}{\partial Y_{22}} \end{aligned}$$

$$= \begin{pmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \\ X_{31} & X_{32} & X_{33} \end{pmatrix} * \begin{pmatrix} \frac{\partial C}{\partial Y_{11}} \\ \frac{\partial C}{\partial Y_{12}} \\ \frac{\partial C}{\partial Y_{21}} \\ \frac{\partial C}{\partial Y_{22}} \end{pmatrix}$$

$$\frac{\partial C}{\partial K} = X * \frac{\partial C}{\partial Y}$$

$$\left\{ \begin{array}{l} Y_1 = B_1 + X_1 * K_{11} \dots X_n * K_{1n} \\ Y_2 = B_2 + X_1 * K_{21} \dots X_n * K_{2n} \\ \vdots \\ Y_d = B_d + X_1 * K_{d1} \dots X_n * K_{dn} \end{array} \right.$$

If we want to calculate $\frac{\partial C}{\partial K_{21}}$ w.g., it only shows up once

$$\frac{\partial C}{\partial K_{21}} = X_1 * \frac{\partial C}{\partial Y_2} \quad (\text{from last example})$$

$$\frac{\partial C}{\partial B_i}$$

$$\frac{\partial C}{\partial b_{ii}} = \frac{\partial C}{\partial y_{ii}} \quad (\text{only appears once})$$

Expanding back,

$$Y_1 = B_1 \dots$$

$$Y_2 = B_2 \dots$$

$$Y_3 = B_3 \dots$$

$$\frac{\partial C}{\partial Y_i} = \frac{\partial C}{\partial B_i}$$

$$\frac{\partial C}{\partial x_j} \Rightarrow Y_i = B_{ij} + x_j + x_{i+1}$$

some $\frac{\partial C}{\partial x_{ij}}$ have different # of terms

\Rightarrow Full cross correlation \Rightarrow Kernel has been rotated 180°

$$\frac{\partial C}{\partial \gamma} = \frac{\partial C}{\partial Y} * k$$

Now, for all γ_i ,

$$\frac{\partial C}{\partial \gamma_i} = \frac{\partial C}{\partial Y_1} * k_{1i} + \dots + \frac{\partial C}{\partial Y_d} * k_{di}$$

$$\frac{\partial C}{\partial \gamma_j} = \sum_{i=1}^d \frac{\partial C}{\partial Y_i} * k_{ij}$$

① pre process the images

② Tokenize the words

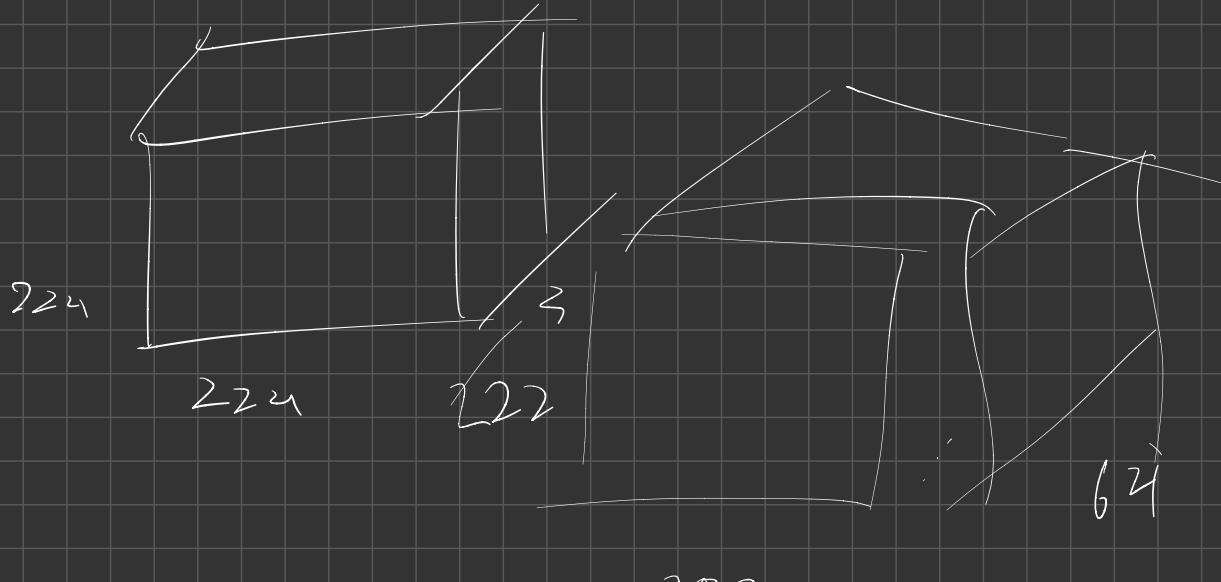
output of CNN \rightarrow feature vector (fixed length)

RNN-based model \rightarrow LSTMs

Initialize hidden state of the RNN with feature vector

Back propagation in time

BLEU METEOR CIDEER, ROUGE



202

