UNIVERSITY OF COPENHAGEN



Today: deep neural networks, convolutional nets



UNIVERSITY OF COPENHAGEN

SILUE .

Julius B. Kirkegaard 2022

The Rectified Linear Unit Activation Function



The Rectified Linear Unit Activation Function

What complex idea could this name possibly encompass?



The Rectified Linear Unit Activation Function

What complex idea could this name possibly encompass?

$$f(x) = \begin{cases} 0, & x \le 0\\ x, & x > 0 \end{cases}$$

















where e.g. f(x) = tanh(x).



- A simple neural network
- It's clear why it's called a *neural* network





- A simple neural network

- It's clear why it's called a *neural* network

Slightly more complex:

$$f(x; p) = \mathbf{a}_1 g(\mathbf{a}_2 g(\mathbf{a}_3 x + \mathbf{b}_3) + \mathbf{b}_2)) + \mathbf{b}_1$$

Some matrix
Some non-linear function
"Activation Function"





- A simple neural network

- It's clear why it's called a *neural* network

Slightly more complex:

$$f(x; p) = \mathbf{a}_1 g(\mathbf{a}_2 g(\mathbf{a}_3 x + \mathbf{b}_3) + \mathbf{b}_2)) + \mathbf{b}_1$$

Some matrix
Some non-linear function
"Activation Function"
Some made much more complex...



Neural networks

In fact, the term "neural network" nowadays can be used about almost any function:





Neural networks

In fact, the term "neural network" nowadays can be used about almost any function:





Neural networks is just ("glorified") function fitting!



Neural networks

Neural networks is just ("glorified") function fitting!

- The main difference is that functions chosen for machine learning application ("neural networks") have **many** parameters that can be tuned ("trained")



Neural networks

If neural network is just function fitting, what do I need to learn?

- Architectures that work well (designing the function) <----- The core of these lectures
- How to fit/train them
- How to train them fast, with lots of data
- How to validate that you've actually "learned" what you think
- How do understand the "reasoning" behind predictions



Deep learning is simply an expression used to indicate that the architecture of our networks are "*deep*".

This is an *architecture* choice that happens to often be a good idea.



















How to do deep learning



"I can do this implement such a function in numpy!"

"... and I can use scipy to fit the function to data!"

"... so can do deep learning!"



How to do deep learning



"I can do this implement such a function in numpy!"

"... and I can use scipy to fit the function to data!"

"... so can do deep learning!"

How to do deep learning



"I can do this implement such a function in numpy!"

"... and I can use scipy to fit the function to data!"

"... so can do deep learning!"

But there are easier ways! (and tons of tricks!)





Ease of use for NNs

(your opinion may differ!)





Ease of doing non-standard stuff

(your opinion may differ!)



PyTorch

PYTÖRCH

Deep Learning with PyTorch

- GPU acceleration
- Automatic Error-Backpropagation (chain rule through operations)
- Tons of functionality built-in





Deep Learning with PyTorch



PYTÖRCH

Deep Learning with PyTorch

What does PyTorch (or other frameworks) do for you?)



Framework feature #1: Automatic gradient calculations

Minimize
$$Loss(f(x,p),y)$$



Framework feature #1: Automatic gradient calculations

Minimize
$$Loss(f(x,p),y)$$

Only feasible way to minimize this if there are many parameters, is if can calculate

$$\nabla D_p Loss(f(x,p),y)$$



Minimize Loss(f(x,p),y)

Loss

Some Paramter p

Minimize Loss(f(x,p),y)





Minimize Loss(f(x,p),y)





Minimize Loss(f(x,p),y)





Minimize Loss(f(x,p),y)



Some Paramter p

Minimize Loss(f(x,p),y)





```
Minimize Loss(f(x,p),y)
```




Requirement 1: Calculate gradients







Requirement 1: Calculate gradients



UNIVERSITY OF COPENHAGEN

This works for basically any function (and it's not magic!)

Framework feature #2: Use GPUs

Graphical processing units are much faster for deep learning than using CPUs!



Framework feature #2: Use GPUs

Graphical processing units are much faster for deep learning than using CPUs!



Requirement 2: GPU

The only change



Requirements for DNN Frameworks

f(x;p)

- Optimisation of parameters p
 - Take first order derivatives
 - Chain rule (backpropagation)
- Process large amounts of data fast
 - Exploit GPUs
- Nice to haves:
 - Standard functions and operations built-in
 - Built-in optimizers
 - Spread training across network
 - Compile for fast inference
 - ...

РҮТ <mark></mark>КСН



Simple PyTorch Example



1	import torch
	V = torch.randn((5, 2))
	W = torch.randn((1, 5))
	def net(x):
	xp = torch.tanh(V @ x)
	🖕 return W @ xp
	x = torch.tensor([0.1, 1.3])
	net(x)



Simple PyTorch Example





More control...



import torch from torch import nn class Net(nn.Module): def __init__(self): super().__init__() self.linear_1 = nn.Linear(2, 5) $self.linear_2 = nn.Linear(5, 1)$ def forward(self, x): xp = torch.tanh(self.linear_1(x)) return self.linear_2(xp)

net = Net()
x = torch.tensor([0.1, 1.3])
net(x)





Χ'

W X'

Some Paramter p

0.55





UNIVERSITY OF COPENHAGEN

f(**V** x)

How to train in PyTorch







How to train in PyTorch



UNIVERSITY OF COPENHAGEN

import torch from torch import nn net = nn.Sequential(nn.Linear(2, 5), nn.Tanh(), nn.Linear(5, 1)) x_data = torch.randn((1000, 2)) y_data = torch.sum(x_data**2, dim=1) iterations = 250 opt = torch.optim.SGD(net.parameters(), lr) for _ in range(iterations): prediction = net(x_data) loss = torch.mean((prediction - y_data)**2) loss.backward() # calculate gradients opt.step() opt.zero_grad()



Better optimiser stepping

$$p_{n+1} = p_n - \nabla_{p_n} \mathcal{L}$$
 · learning rate

- What if some gradients are much smaller than others?
- What happens when gradients disappear when loss is small?



Better optimiser stepping

$$p_{n+1} = p_n - \nabla_{p_n} \mathcal{L}$$
 · learning rate

- What if some gradients are much smaller than others?
- What happens when gradients disappear when loss is small?

Solution : Variable learning rates and momentum

• Many algorithms exists, perhaps most popular: "Adam"



Better optimiser stepping

SGD (Stochastic gradient descent)





Adam (Adaptive Moment Estimation)



	opt	= torch.optim.Adam(net.parameters(), lr)
	for	_ in range(iterations):
		prediction = net(x_data)
		<pre>loss = torch.mean((prediction - y_data)**2)</pre>
		loss.backward()
		opt.step()
		opt.zero_grad()



A more complex example:

Neural networks are universal function approximators





Hidden Layer $\in \mathbb{R}^5$ Output Layer $\in \mathbb{R}^2$



Broad network



Loss functions

Regression

"Output of network is a number that should be close to label y

L2 loss:

$$L(x, y) = (NN(x) - y)^2$$



Loss functions

Regression

"Output of network is a number that should be close to label y

L2 loss:

$$L(x,y) = (NN(x) - y)^2$$

Classification

"Output of network is a probability distribution over labels ($y_i = 1$ for true label and = 0 otherwise).

Cross-entropy loss:

$$\mathbf{p} = \mathrm{NN}(x)$$
$$L(\mathbf{p}, \mathbf{y}) = -\sum_{i} y_i \, \log(p_i)$$





Exemplified using image classification..





Find images with angle:







Find images with angle:







Find images with angle:



We don't care where in the images the line is! (like: *"is it a cat or not?"*)





What we want:



-) = large number (high probability)
 - NN(



) = large number (high probability)



) = small number (low probability)

NN(

-) =
 -) = small number (low probability)



What we want:



How can we design such a function?









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





Place kernel somewhere on image, multiply and sum:



= small number





Place kernel somewhere on image, multiply and sum:



= small number



= large number





What happens for a "wrong" image?



= small number



= small(ish) number



Idea: try all locations of kernel and find maximum:



Idea: try all locations of kernel and find maximum:





Idea: try all locations of kernel and find maximum:





Idea: try all locations of kernel and find maximum:





This is precisely what a convolution does!





This is precisely what a convolution does!




But we didn't train anything?





But we didn't train anything?







Run code!

But we didn't train anything?







(slightly smarter than ours!, albeit noisy)



But we didn't train anything?



We then make it deep, so kernels can be combined:





But we didn't train anything?



We then make it deep, so kernels can be combined:



The number "4" is made up of four small lines



But we didn't train anything?



We then make it deep, so kernels can be combined:

For deep networks, this max is taken over a few neighbors at a time, not entire image.



The number "4" is made up of four small lines









https://poloclub.github.io/cnn-explainer/ http://www.cs.cmu.edu/~aharley/nn_vis/cnn/2d.html



CNN: The Building Blocks



CNN: The Building Blocks





CNN: The Building Blocks





CNN: A *deep* example:

```
nn.Sequential(nn.Conv2d(1, 32, 7),
  nn.BatchNorm2d(32),
  nn.ReLU(),
  nn.MaxPool2d(3),
  nn.Conv2d(32, 32, 3),
  nn.BatchNorm2d(32),
  nn.ReLU(),
  nn.Conv2d(32, 3, 3),
  nn.MaxPool2d(3),
  nn.Flatten(),
  nn.Linear(432, 20),
  nn.ReLU(),
  nn.Linear(20, 1),
  nn.Sigmoid())
```

