

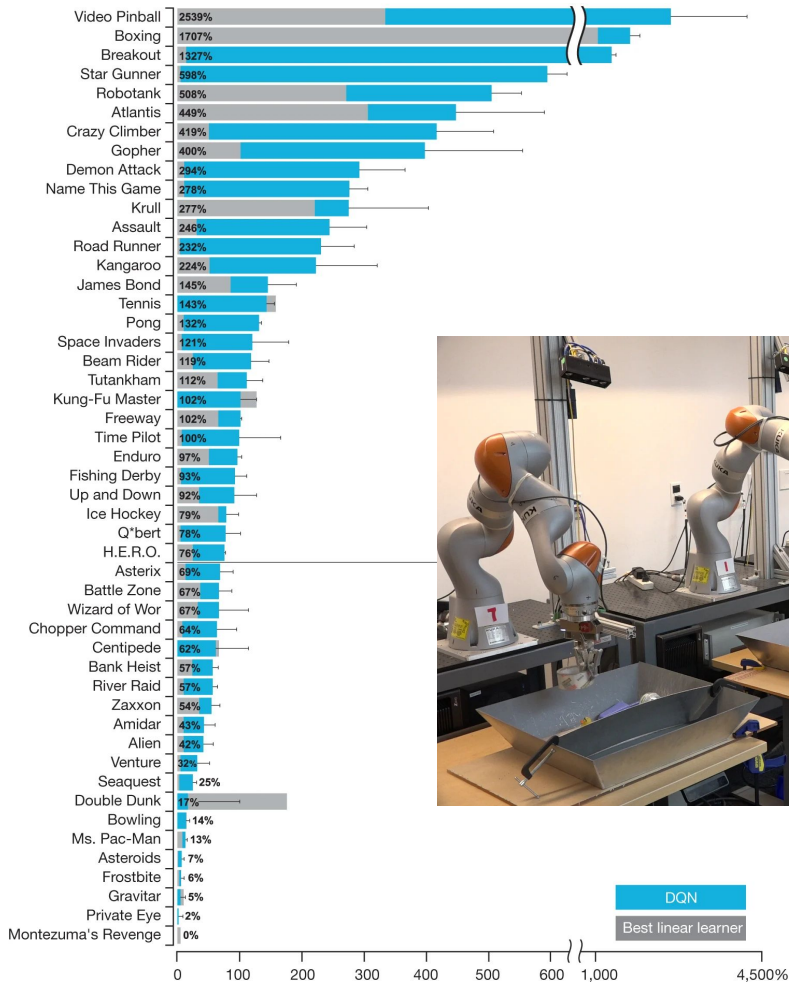
Deep Learning and PyTorch Tutorial

17th Jan 2023
Skanda Vaidyanath
CS 234

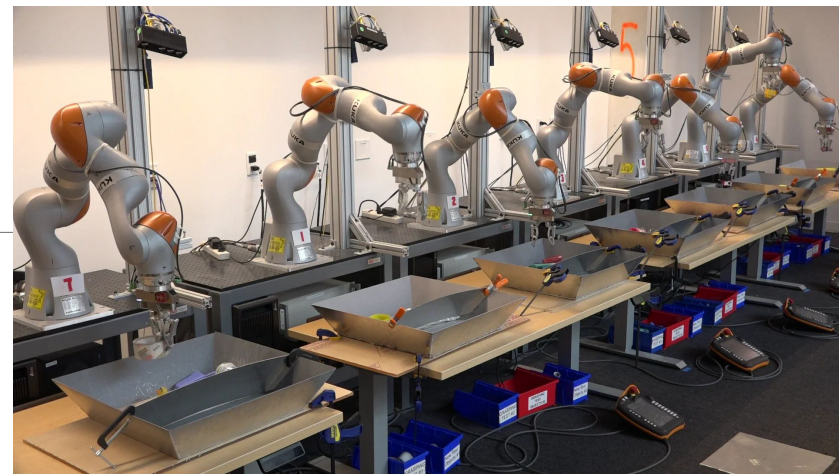
Overview

- Motivation
- Neural networks
- Backpropagation
- Training and Optimization
- CNNs
- PyTorch

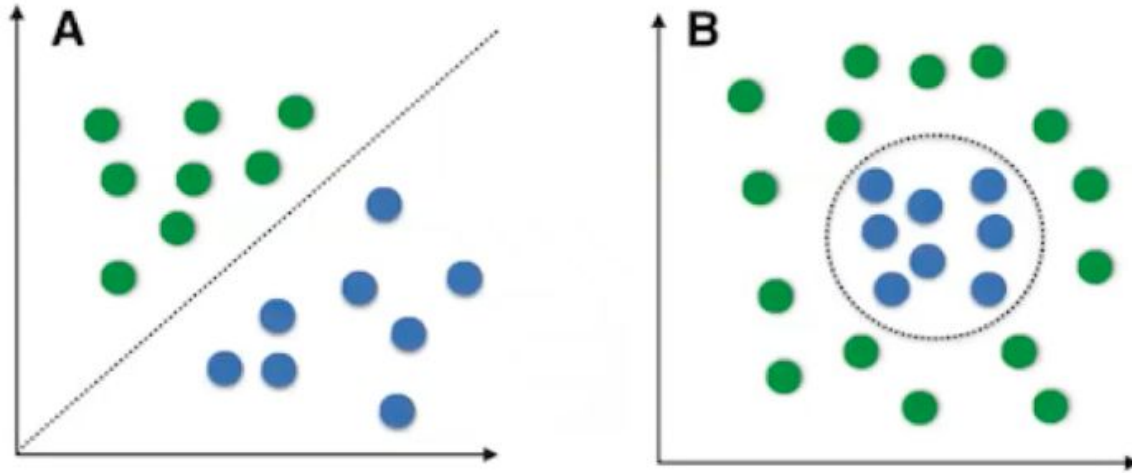
Motivation



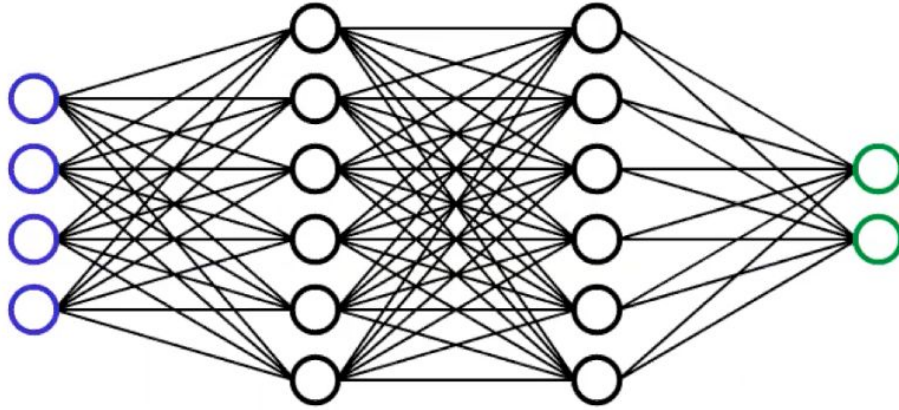
[Image source](#)



The need for non-linear models



Neural networks



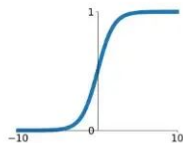
$$x \qquad h_1 = \phi(W_1x + b_1) \qquad h_2 = \phi(W_2h_1 + b_2) \qquad y = \phi_{out}(W_3h_2 + b_3)$$

Non-linearities

Activation Functions

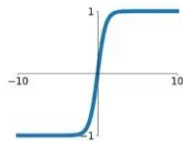
Sigmoid

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



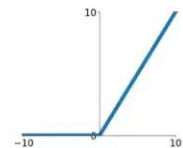
tanh

$$\tanh(x)$$



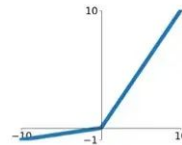
ReLU

$$\max(0, x)$$



Leaky ReLU

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



Maxout

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

ELU

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

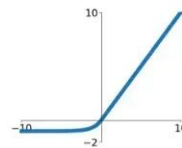


Image credit: [Medium article](#)

Computational graphs

Eg: $f(x,y,z) = (x+y)z$

What is the derivative of f wrt x, y, z ?

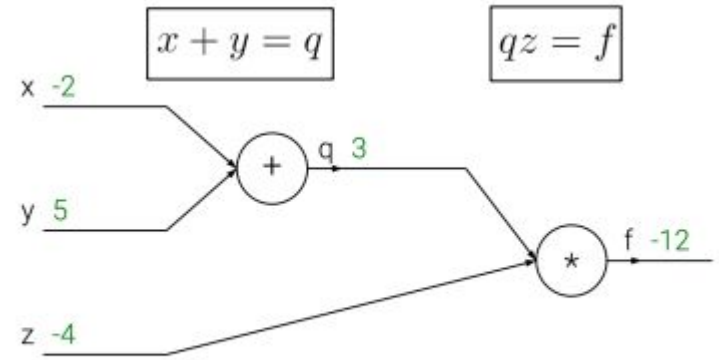
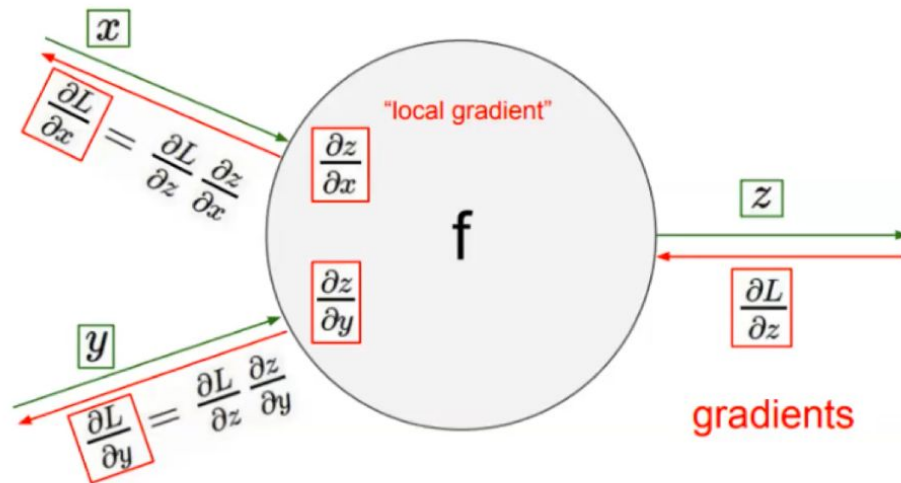


Image credit: [CS231N Slides](#)

Backpropagation

Using the chain rule to calculate gradients



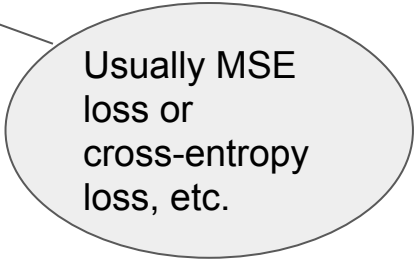
Training a neural network

Loss function, usually an average over data points:

$$L(\theta) = \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \ell(\theta; x, y)$$

Want to find

$$\theta^* \in \arg \min_{\theta} L(\theta)$$



Usually MSE
loss or
cross-entropy
loss, etc.

Optimization

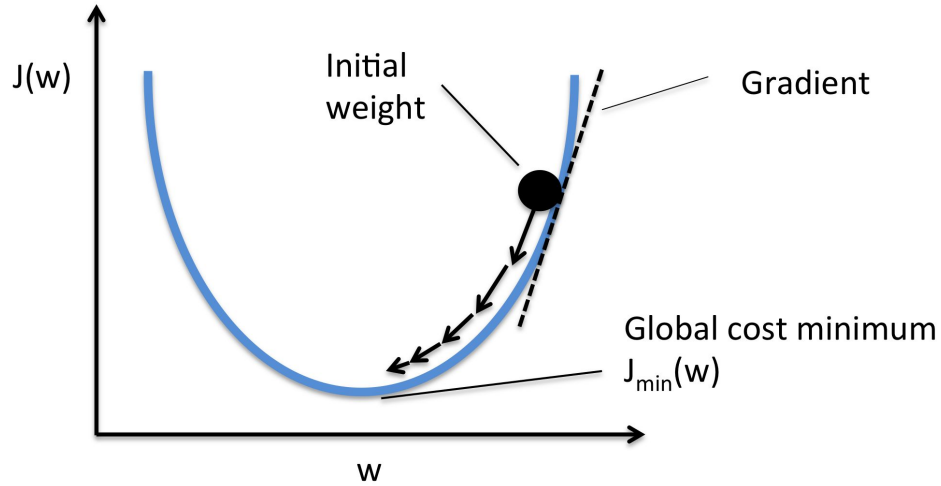


Image credit: [Blog post](#)

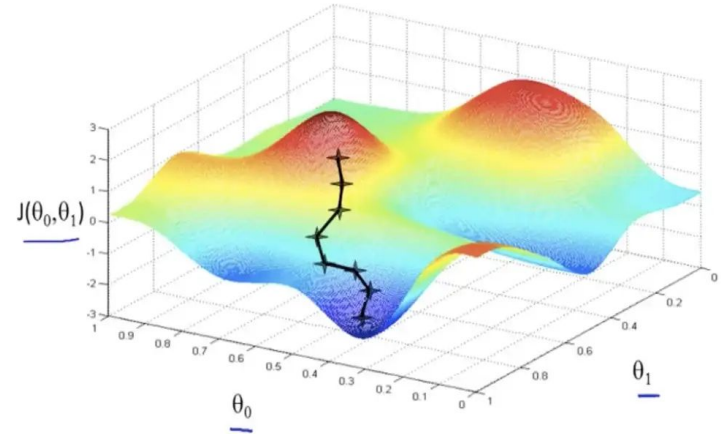


Image credit: [Medium article](#)

Stochastic Gradient Descent

Often it is infeasible/ inefficient to do full batch gradient descent, so we resort to stochastic gradient descent (SGD)

1. Pick some θ_0 , then iterate:
2. $\theta_{k+1} = \theta_k - \alpha_k \nabla L(\theta_k)$

$$\nabla L(\theta) = \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \nabla_{\theta} \ell(\theta; x, y) \quad O(|\mathcal{D}|)$$

SGD: sample mini-batch $B \subseteq \mathcal{D}$, use

$$\hat{\nabla} L(\theta) = \frac{1}{|B|} \sum_{(x,y) \in B} \nabla_{\theta} \ell(\theta; x, y) \quad O(|B|)$$

Note $\mathbb{E}[\hat{\nabla} L(\theta)] = \nabla L(\theta)$

Other optimizers

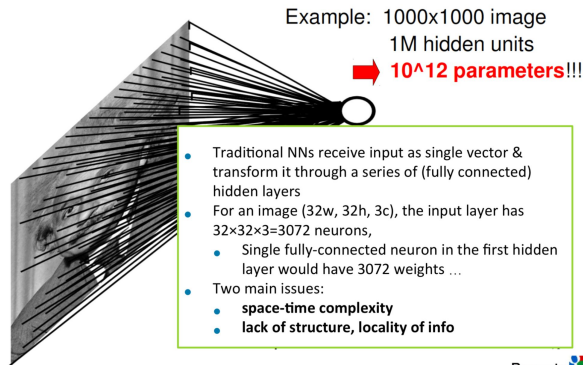
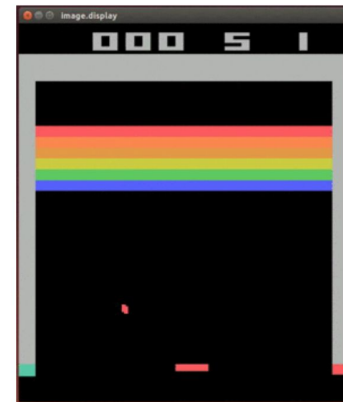
There are several new optimizers that improve on vanilla gradient descent. Some of them are:

- Momentum
- Adagrad
- RMSProp
- Adam

To learn more: <https://ruder.io/optimizing-gradient-descent/>

Convolutional Neural Networks

- CNNs are extensively used in computer vision
- If we want to go from pixels to decisions, likely useful to leverage insights from visual input
- Why not just use fully connected networks?



Convolutional Neural Networks

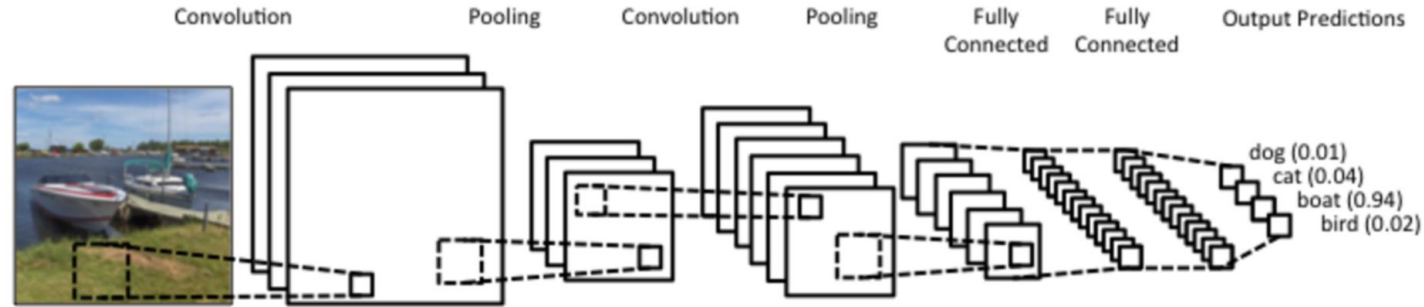
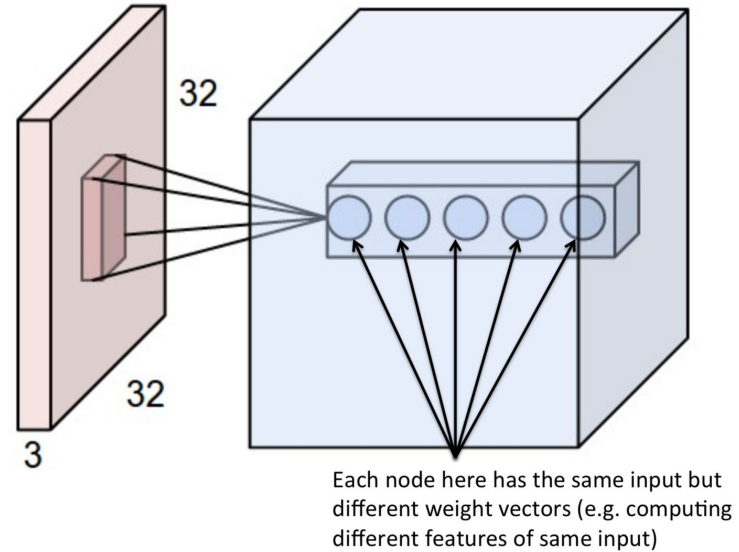


Image: <http://d3kbpzbmcyntmx.cloudfront.net/wp-content/uploads/2015/11/Screen-Shot-2015-11-07-at-7.26.20-AM.png>

Convolutional Neural Networks

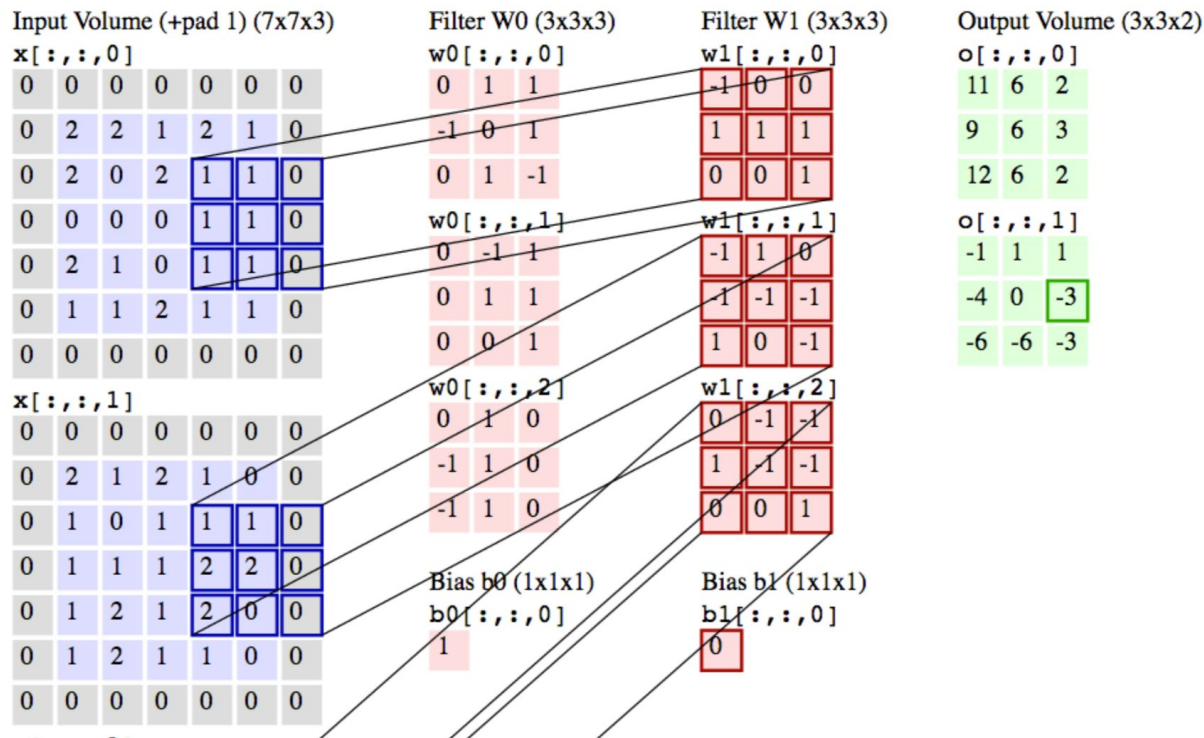
- Image dimensions
- Filters/ kernels
- Stride
- Padding



Convolutional Neural Networks

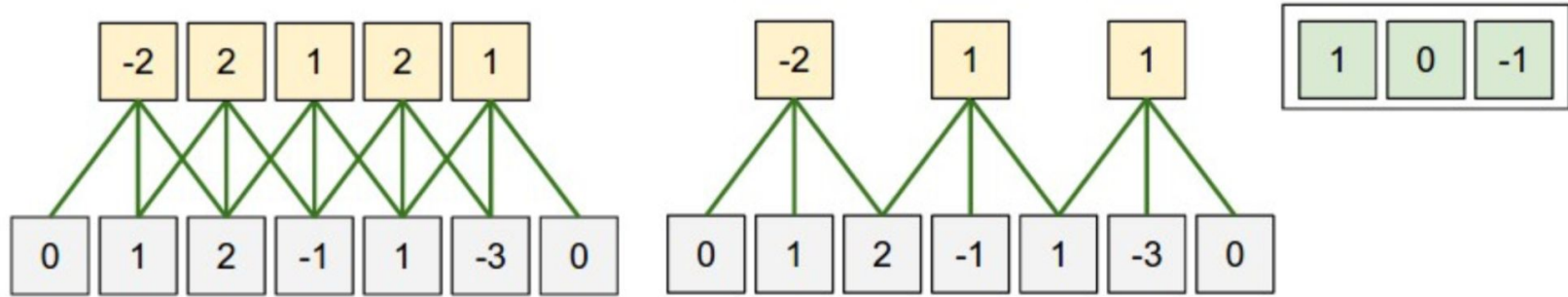
Full demo:

<https://cs231n.github.io/convolutional-networks/>



Convolutional Neural Networks

Simpler example:



PyTorch



- Objective: introduce key PyTorch concepts for CS234
- PyTorch documentation: <https://pytorch.org/docs/stable/index.html>
- Tutorial: Python notebook will be shared after class