

Lecture 1 (History, Evolution and Revolution of DL) 2/12/2020

Inspiration of DL: The Brain

- McCulloch & Pitts (1943): networks of binary neuron can do logic
- Donald Hebb (1947): Hebbian synaptic plasticity
- Norbert Wiener (1948): cybernetics, optimal filter, feedback, auto-poiesis, auto-organization
- Frank Rosenblatt (1957): perceptron
- Hubel & Wiesel (1960s): visual cortex architecture

stopped and again started in 1980s and died in 1995 again - and again rose up in 2010s with speech recognition.

2012-2013 → rise of Computer Vision

2016 → rise of NLP

Supervised Learning: Training a machine by showing examples instead of programming it.

Works well for:

Speech → words

Image → categories

Portrait → name

Photo → caption

Text → topic

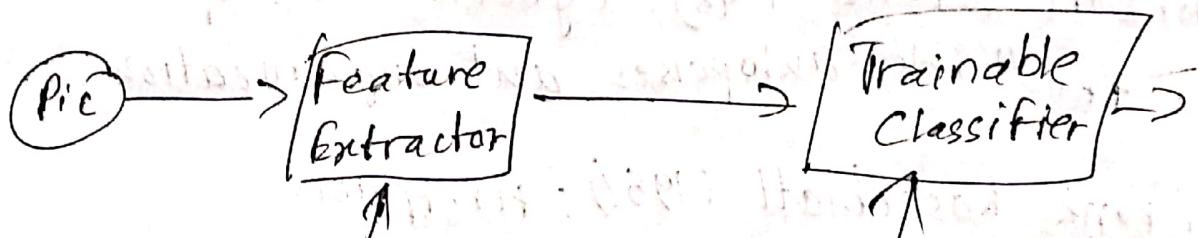
Supervised Learning goes back to the Perceptron and Adaline

The McCulloch-Pitts Binary Neuron

→ Perceptron: Weights are motorized potentiometers

→ Adaline: Weights are electrochemical memistors

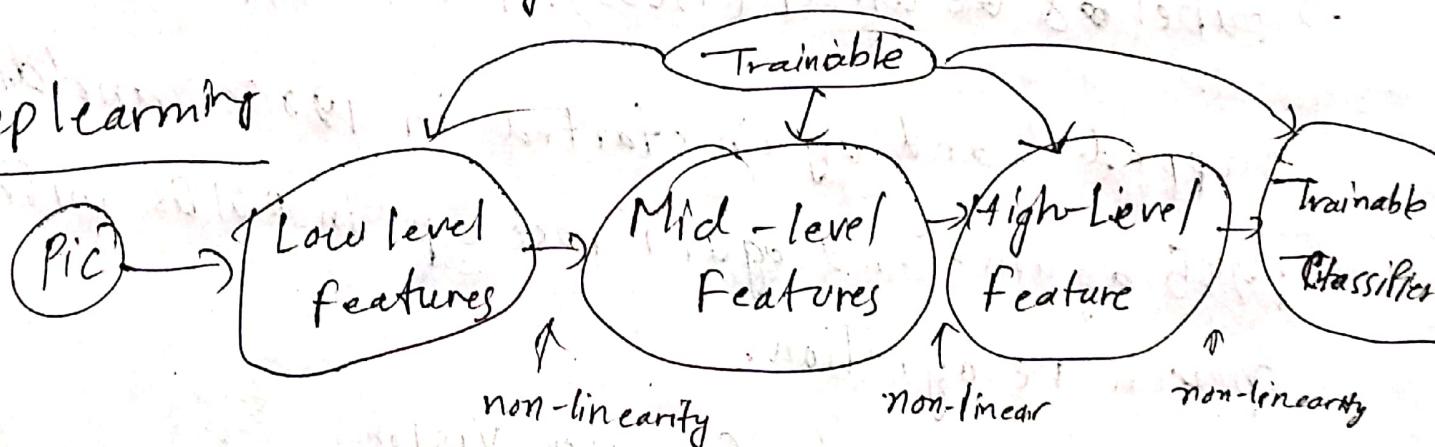
Traditional Machine Learning



Hand engineered.

Trainable

Deep learning



MLP (Multi-layer Perceptron)

$$\text{ReLU}(x) = \max(x, 0)$$

Supervised Machine learning = function optimization

Stochastic gradient descent

$$w_i \leftarrow w_i - \eta \frac{\partial L(w, X)}{\partial w_i}$$

Computing Gradients by Back-Propagation:

$$\frac{dC}{dW_n} \longrightarrow$$

Traditional NN has usability problems and many others dealing with images.

- Hubel and Wiesel's Model of the Architecture of the Visual System
 - simple cells detect local features [1962]
 - complex cells "pool" the outputs of simple cells within space
↳ They found out how things work in brains for recognition

[Fukushima 1982]

[LeCun 1989, 1998]

ConvNets can recognize multiple objects

- All layers are convolutional
- Networks performs simultaneous segmentation and recognition

Face and Pedestrian Detection with ConvNets (1993-2015)

Training a Robot to Drive Itself in Nature [Hadsell 2009]

Semantic Segmentation with ConvNets [Farabet 2012]

1986-1996] \Rightarrow Special NN chips

Deep learning Revolution

Speech Recognition 2010

Image Recognition 2013

NLP 2015

Deep ConvNets for Object Recognition (on GPU)

AlexNet [2013] OverFeat [2013]

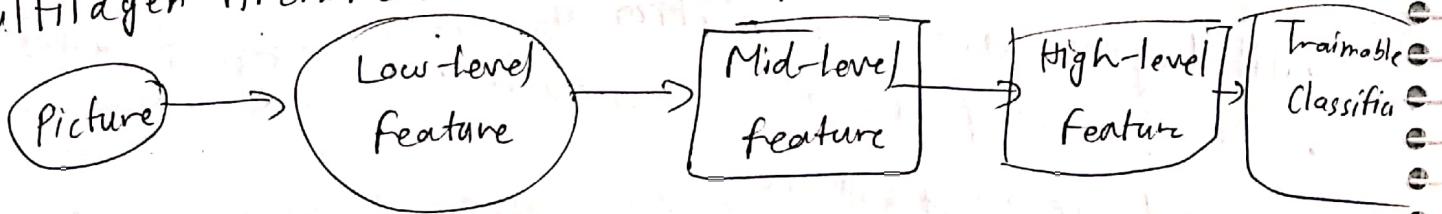
VGG [2013]

GoogleNet [2013]

ResNet [2015]

DenseNet [2017]

Multilayer Architectures == Compositional structure of Data



MSRA - 2015

Mask R-CNN
2017

Mask R-CNN
2017

② Mask-R-CNN on COCO dataset

→ Individual objects are segmented

3D ConvNet for Medical Image Analysis

Deep learning enables

→ Safer cars, autonomous cars

→ Better medical image analysis

→ Personalized medicine

→ Adequate language translation

→ Useful but stupid chatbots

→ Info search, retrieval, filtering

→ other

It can't yet

Machines with common sense

Intelligent personal assistants

Smart chatbots

Household robots

Agile and dexterous robots

Artificial General Intelligence

Intelligence

Deep learning = Learning Representations / Features

Hierarchical representation

→ Hierarchy of representations with increasing level of abstraction

→ Each stage is a kind of trainable feature transform

SVM is nothing but two layer neural nets and
the first layer is trained in an unsupervised way

→ Deep machines are more efficient for representing
certain classes of functions.

Why would DL be more efficient? [Bengio, LeCun 2007]

→ less params needed to get good result

What are Good Features?

Discovering the Hidden Structure in High Dimensional

Data: The Manifold Hypothesis

The Manifold Hypothesis

- Natural data lives in a low-dimensional manifold [non-linear]
- Because variables in natural data are mutually dependent

Disentangling factors of variation

- The ideal disentangling feature Extractor
- , PCA can find the representation if they are linear

Practicum 1

[Yann LeCun, Alfredo Canziani]
Mark Goldstein

Classification, linear algebra and visualization

We have to find the point of classification to and bring it to the main focus.

Matrix multiplication: (what does it do?) [Linear transformation]

- Rotation
 - Stretching/Scaling (Zooming / shearing)
 - Reflection
 - Translation
 - Affine transformation
- Neural net does this by trying to reach convergence

• We need non-linearity because without non-linearity we can't get

Space Stretching

We aim to divide the space so that we can get linear separable translation

Week 21 Lecture: SGD and Backprop

Parameterized Model

$$\hat{y} = G(x, w)$$

params weights

[Parameterized Deterministic Function]

Computing function G may involve complicated algo.

Variables (tensor, scalar, continuous, discrete, ...)

Deterministic function

Scalar valued function (implicit output)

Loss function and average loss

, machine learning is all about optimizing functions.

Sto. gradient descent

$$w \leftarrow w - \eta \frac{\partial L(s, w)}{\partial w}$$

SGD exploits the redundancy in the samples

. we use minibatch for parallelization in practice

There are optimization algorithms that are not dependent on gradient specially the ones where we don't know the grads or can't get the grads. Those are called zeroth order method.

- In RL, the output function is not differentiable most of the time, which makes it really inefficient.
- Reward is the negative of cost. in RL. RL practitioners generally use a secondary cost function which is differentiable to make RL efficient.

Traditional NN

stacked linear and non-linear functional blocks

$$s[i] = \sum_{j \in \text{UP}(i)} w[i,j] z[j] \quad z[i] = f(s[i])$$

Backprop through a non-linear functions

$$\text{Chainrule } g'(h(s))' = g'(h(s)) \cdot h'(s)$$

Perturbations

Perturbations:

Linear blocks $s_{k+1} = w_k z_k$

Nonlinear blocks $z_k = h(s_k)$

PyTorch definition: (object oriented version)

nn.Linear automatically adds a bias factor

Jacobian Matrix

Partial derivatives of i th output wrt j th input

$$\left(\frac{\partial z_g}{\partial z_f} \right)_{ij} = \frac{(\partial z_g)_i}{(\partial z_f)_j}$$

Backprop through a multi-stage graph

Two Jacobian matrices

one with respect to $z[k] \rightarrow$ input

One with respect to $w[k] \rightarrow$ weight

Log softmax $y_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$ [vector of probabilities] $\log y_i$
John Vidal

$$\log y_i = x_i - \log \sum_j e^{x_j}$$

Practical tips of Backprop

softmax is just a generalization of
soft sigmoid with multiple outputs

- Use ReLU . (tanh and others are pd falling out of favor)
- Use cross entropy for classification
- Use sgd on minibatched
- Shuffle the training samples
- Normalize the input variables (zero mean, unit variance)
- Schedule to decrease the learning rate
- Use a bit of L1 ^{or} L2 regularization on the weights
(Weight Decay)
→ it's best to turn it on after couple epochs
- Use dropout for regularization (distribute info better)
- Lecun → Efficient Backprop (1998)
- Neural Nets, Tricks to trade
 - ~~Weight Decay~~
- Any directed acyclic graph is ok for backprop

Practicum 2: Training a neural network

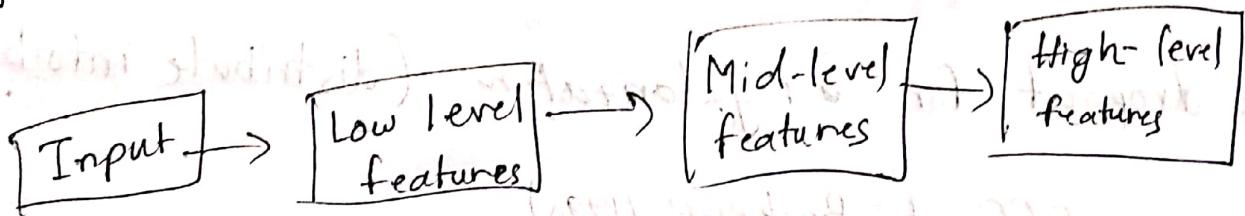
NN are mainly doing rotation and squashing

ANN-supervised learning

$$x_c(t) = t \begin{pmatrix} \sin\left(\frac{2\pi}{C}(2t+c-1)\right) \\ \cos\left(\frac{2\pi}{C}(2t+c-1)\right) \end{pmatrix} + N(0, \sigma^2)$$

$$0 \leq t \leq 1, \quad c=1, \dots, C$$

we try to adapt the decision boundaries by rotating and scaling



If we have multiple output then ~~it's~~ it is called multi-headed network

Train data

$$x^{(1)} \in \mathbb{R}^n$$

$$x = \begin{bmatrix} -x^{(1)} \\ -x^{(2)} \\ \vdots \\ -x^{(m)} \end{bmatrix} \quad m$$

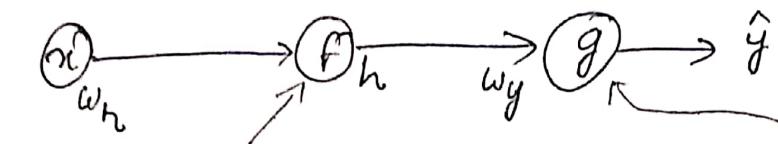
$$c_i \in \{1, 2, \dots, K\}$$

$$C = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_m \end{bmatrix} \quad m$$

$$\text{1-hot encoding} \quad k$$

$$y = \begin{bmatrix} (1, 0, 0) \\ (0, 1, 0) \\ \vdots \\ (0, 0, 1) \end{bmatrix} \quad m$$

Fully connected layer



$$h = f(w_h x + b_h)$$

$$\hat{y} = g(w_y h + b_y)$$

Affine transformation \rightarrow rotation and scaling
Non-linear function \rightarrow squashing

$$w_h \in \mathbb{R}^{d \times n}$$

$$f, g = (\cdot)^+, \delta(\cdot)$$

$$b, b_h \in \mathbb{R}^d$$

$$\tanh(\cdot), \text{softmax}, \text{max}(t)$$

$$w_y \in \mathbb{R}^{K \times d}$$

- Backward pass computes the grads of loss w.r.t. our learnable params

$$\text{softmax}(\ell)[c] = \frac{\exp(\ell[c])}{\sum_{k=1}^K \exp(\ell[k])} \in (0, 1)$$

$$L(\hat{y}, c) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, c_i)$$

$$L(\hat{y}, c) = -\log(\hat{y}[c])$$

↑ cross entropy / negative log likelihood

$$x, c=1 \Rightarrow y = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$

$$\theta = \{w_h, b_h, w_y, b_y\}$$

$$\hat{y} \left(\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \right) \Rightarrow L \left(\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, 1 \right) \rightarrow 0^+$$

$$\mathcal{J}(\theta) = L(\hat{y}(\theta), c) \in \mathbb{R}^+$$

$$\hat{y}(x) = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \Rightarrow L \left(\begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, 1 \right) \rightarrow +\infty$$

$$-\frac{d\mathcal{J}(\theta)}{dy} (y_0)$$

- notebooks
- spiral classification
- regression

choice of activation will have effects on regression

b)

Lecture 3: Convolution Neural Networks

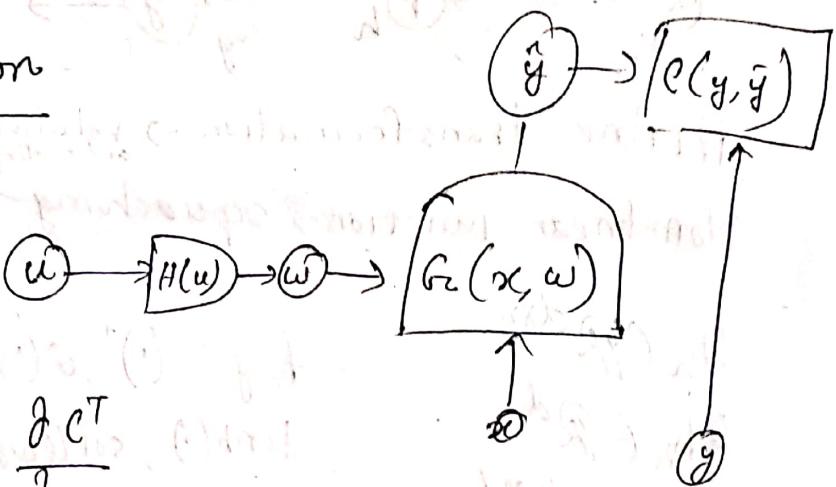
Parameters transformation

$$w = H(u)$$

Parameter vector output of a function can be the

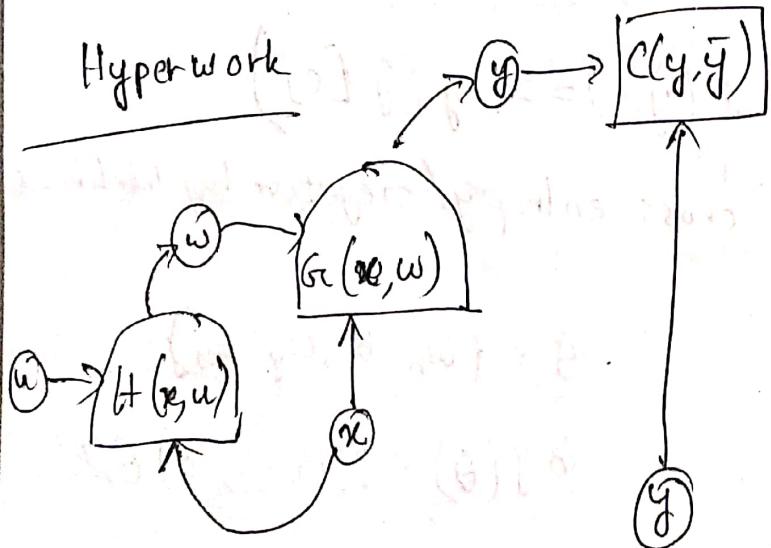
$$u \leftarrow u - \eta \frac{\partial H^T}{\partial u} \frac{\partial C^T}{\partial w}$$

$$w \leftarrow w - \eta \frac{\partial H}{\partial u} \frac{\partial H^T}{\partial u} \frac{\partial C^T}{\partial w}$$



$$[N_w \times N_u] \quad [N_u \times N_w] \quad [N_w \times 1]$$

Hyperwork



Shared weights

These can be necessary for motif detection anywhere on an input

• When the parameter vector in the output of another network $H(x, u)$

• The weights of network $G_r(x, w)$ are dynamically configured by network $H(x, u)$

Conv Basics

Stride, pad

$$y_i = \sum_j w_j x_{i,j} \quad (\text{conv})$$

$$y_i = \sum_j w_j x_{i,j} \quad (\text{cross correlation})$$

Deep learning - Learning Hierarchical Representations

Idea of convnets comes from biology.

(How does the brain interprets images?)

• Neuron firing

Hubel and Wiesel 1962

• simple cells detect local features

• complex cells "pool" the output of simple

First 'Real' Convnet at Bell Labs [Lecun 1989]

Conv + Pool
ReLU

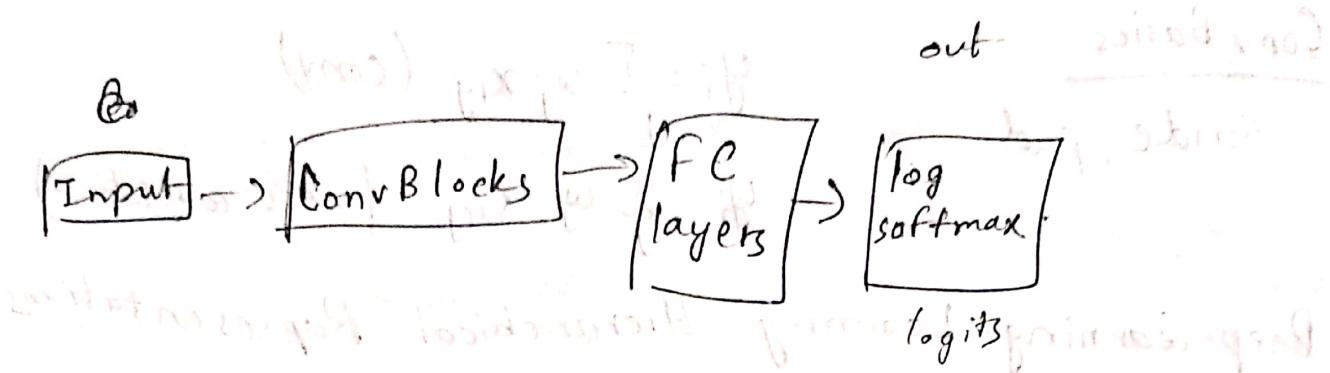
Normalization
Batch-Norm

Filter Banks
Conv

Non-linearity
ReLU

Pooling

ConvNet Block



LeNet (@Character recog)

Multiple Character recognition [Malan 1992]

features to get detect something

Feature Binding and What convnets are good for

- Signals that come to you in the form of images
- Signals that have strong local correlations
- Signals where features can appear anywhere
- Signals in which objects are invariant to translations and distortions
- 1D convnets: sequential signals, text
 - Text, music, audio, speech, time series
- 2D convnets: images, fine frequency representation (audio & speech)
 - Object detection, localization, recognition
- 3D convnets: video, volumetric images, tomography images
 - Video recognition/understanding, hyperspectral image analysis
 - Bio image analysis

Practicum 3 Natural signals properties and CNNs

Foundation of DL

CNN (exploiting stationarity, locality and compositionality of natural data)

$$x = [x_1 \ x_2 \ x_3 \dots \ x_t \dots]^T \quad [\text{Time series}]$$

x_t are waveform heights

$$x = [x_{11} \ x_{12} \ \dots \ x_{1n} \ x_{21} \ x_{22} \dots]^T \quad [\text{Image}]$$

x_{ij} are pixel values

$$x = [x_1 \ x_2 \ x_3 \dots \ x_n]^T \quad [\text{Text}]$$

x_t are one hot vectors

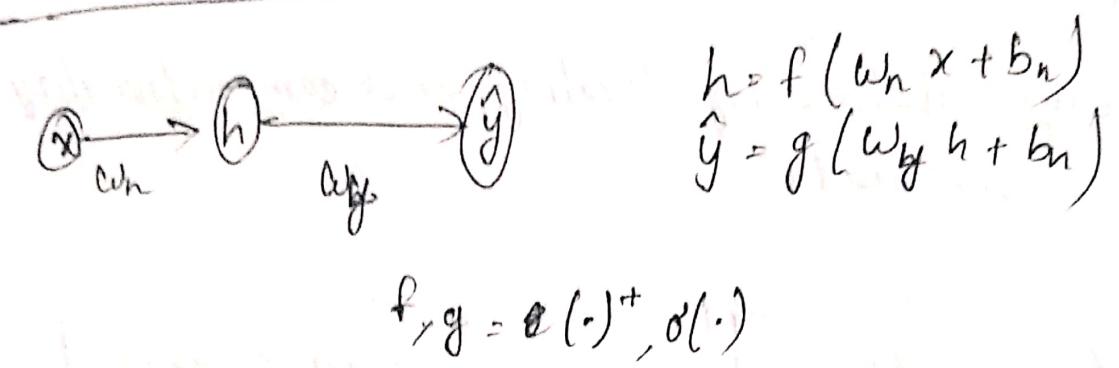
3 properties of natural data

Stationarity \rightarrow we observe the same pattern over and over in our data

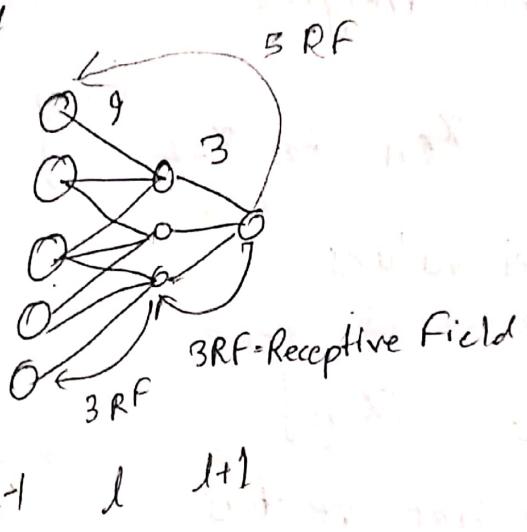
Locality \rightarrow the patterns are localized

Compositionality \rightarrow perceiving things with composing

FC(fully connected) Layers

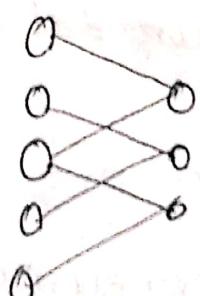
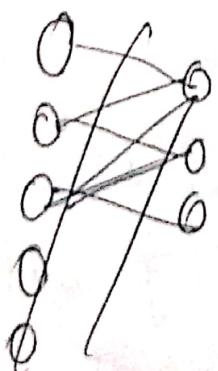


Locality \Rightarrow sparsity



- reduced amount of computation

Stationary \Rightarrow parameters sharing



- faster convergence
- better generalization
- constrained to input size
- kernel independence
 \Rightarrow high parallelizability

Kernels - 1D data

Padding 1D data

Standard spatial CNN

• Multiple layers

• Convolutions

• Non-linearity (ReLU and Leaky)

• Pooling (Geoffrey Hinton don't like this)

• Batch normalisation

• Residual bypass connection

He is working on it
to make this more

so that we don't
need to ~~through~~
away information

• If we have the 3 properties we can use CNN

Or else it's not appropriate

Practicum 4: Listening to convolution [1D convs]

Linear Algebra Review

$$\underline{h} = f(\underline{x})$$

$$h = f(x)$$

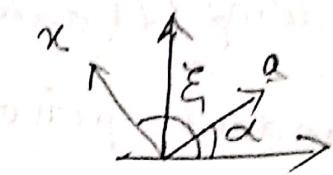
$$\underline{z} = \underline{A} \underline{x} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} =$$

$$= \begin{bmatrix} -\underline{\alpha^{(1)}} \\ -\underline{\alpha^{(2)}} \\ \vdots \\ -\underline{\alpha^{(m)}} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

$$= \begin{bmatrix} \underline{\alpha^{(1)}x} \\ \underline{\alpha^{(2)}x} \\ \vdots \\ \underline{\alpha^{(m)}x} \end{bmatrix}$$

$\alpha^T x$ when $n=2$

$$\alpha^T x = \alpha_1 x_1 + \alpha_2 x_2$$



$$= \|\alpha\| \cos \alpha \|x\| \cos \xi + \|\alpha\| \sin \alpha \|x\| \sin \xi$$

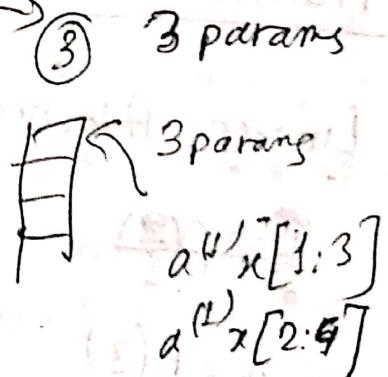
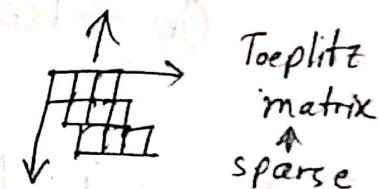
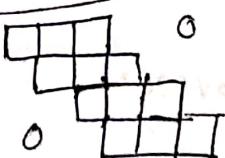
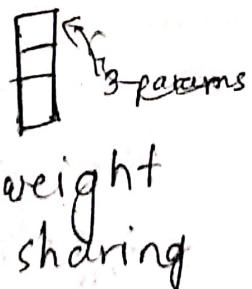
$$= \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_m \end{bmatrix} = \underline{z}$$

$$= \|\alpha\| \|x\| (\cos \alpha \cos \xi + \sin \alpha \sin \xi)$$

$$= \|\alpha\| \|x\| \cos(\xi - \alpha)$$

$$z = \underline{\alpha_1} x_1 + \underline{\alpha_2} x_2 + \dots + \underline{\alpha_n} x_n$$

Extension to convolutions:



Toepplitz
matrix
sparse

Convolution is a matrix multiplication
with a lot of zeros

Listening to convs

Mathematical formalism

$\mathcal{X} = \{x^{(i)} \in \mathbb{R}^n \mid x^{(i)} \text{ is a data sample}\}_{i=1}^m$, input sample

$\mathcal{X} = \{x^{(i)} : \Omega \rightarrow \mathbb{R}^c, \omega \mapsto x^{(i)}(\omega)\}_{i=1}^m$

$\Omega = \{1, 2, \dots; T/\Delta t\} \subset \mathbb{N}$

$\Omega \times \{1, \dots, h\} \times \{1, \dots, w\} \subset \mathbb{N}^2$

$\Omega = \mathbb{R}^4 \times \mathbb{R}^4, c=1$ (Hamiltonian)
space-time four-momentum

Week 5

Lecture 4: Optimization

Aaron Defazio (Facebook Research Scientist)

- Gradient descent & SGD
- Momentum
- Diagonal rescaling: RMSProp & Adam
- Normalization layers
- Solving optimization problems without optimizing
- Application: MRI reconstruction

Gradient descent: stepsize

$$w_{k+1} = w_k - \gamma_k \nabla f(w_k)$$

$$\underset{w}{\text{minimize}} \quad f(w)$$

follow direction of the negative gradient
"worst optimization algo"

The negative gradient direction is the steepest direction locally

SGD:

$$w_{k+1} = w_k - \gamma_k \nabla f_i(w_k)$$

SGD: often the best method available

$$E[\nabla f_i(w_k)] = \nabla f(w_k)$$

SGD is GD with noise

Advantages of SGD:

- There is redundant info across instances the noise is lower when there is higher redundancy
- At early stages of optimization, the noise is small compared to the information in the gradient so a SGD step is virtually as good as a GD step
- The noise can prevent the optimizing converging to bad local minima, a phenomena called annealing
- Stochastic gradients are drastically cheaper to compute (proportional to your dataset size), so you can often take thousands of SGD steps for the cost of one GD steps.

Mini-batching

$$w_{k+1} = w_k - \frac{\gamma}{k|B|} \sum_{j \in B} \nabla f_j(w_k)$$

- Done for parallelization in GPU

Theoretically is not right

Empirically gives a good result

For Batch optimization, LBFGS & Conjugate

Gradients are far better than vanilla SGD.

Momentum:

$$w_{k+1} = w_k - \gamma_k \nabla f_i(w_k) + \beta (w_k - w_{k-1}) \quad \begin{matrix} \text{SGD +} \\ \text{Momentum} \end{matrix}$$

$$p_{k+1} = \hat{\beta}_k p_k + \nabla f_i(w_k)$$

$$w_{k+1} = w_k - \gamma_k p_{k+1}$$

- The next step becomes a combination of the previous steps direction and the new negative gradient

optimization process resembles a heavy ball rolling down a hill.

Larger momentum = slower reaction to change in the landscape

$\beta = 0.9$ or 0.99 almost always works well.

- step size need to be decreased when the momentum parameter is increased to maintain convergence

why momentum works?

More explanation one

Acceleration

Nesterov's momentum

$$p_{k+1} = \hat{p}_k p_k + \nabla f_i(w_k)$$

$$w_{k+1} = w_k - \gamma_k (\nabla f_i(w_k) + \hat{p}_k p_{k+1})$$

→ Nesterov

Explanation two: Noise Smoothing

In optimization, we

$$\bar{w}_k = \frac{1}{K} \sum_{k=1}^K w_k$$

smooths the noise from SGD

when using carefully chosen constants, probably "accelerates" convergence on problems with simple structure
torch.optim.SGD(..., nesterov=True)

Adaptive methods

The magnitude of the gradients often varies highly between layers due to, so a global learning rate may not work well.

- Instead of using the same learning rate for every weight in our network, maintain an estimate of a better rate separately for each weight.
- most methods either adapt to the variance of the weights, or to the local curvature of the problem

RMSProp (Root Mean Square Propagation)

$$V_{t+1} = \alpha V_t + (1-\alpha) f_i(w_t)^2$$

$$w_{t+1} = w_t - \frac{\Delta f_i(w_t)}{\sqrt{V_{t+1}} + \epsilon}$$

- if the expected value of gradient is small, this is similar to dividing by the standard deviation

Adam (Adaptive moment estimation)

RMSProp with a kind of momentum
momentum (exponentially moving average)

$$m_{t+1} = \beta v_t + (1-\beta) \nabla f_i(w_t)$$

Same as
RMSProp

$$v_{t+1} = \alpha v_t + (1-\alpha) \nabla f_i(w_t)^2$$

$$w_{t+1} = w_t - \eta \frac{m_t}{\sqrt{v_{t+1}} + \epsilon}$$

Recommended
SGD + Momentum
Adam

- The exponentially moving-average method of updating momentum is equivalent to the standard form under rescaling.

→ full version of Adam has bias-correction, which just keeps the moving averages unbiased during early iterations

Disadvantages of Adam

i) doesn't converge at all
problems

ii) gives worse generalization on some computer vision problems than SGD

iii) requires more memory than SGD

iv) Need to tune both momentum params

on some simple example

Normalization layers

Learnable scale factor input

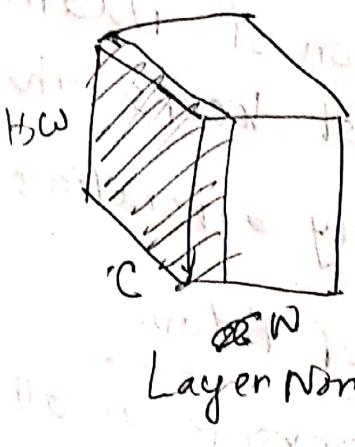
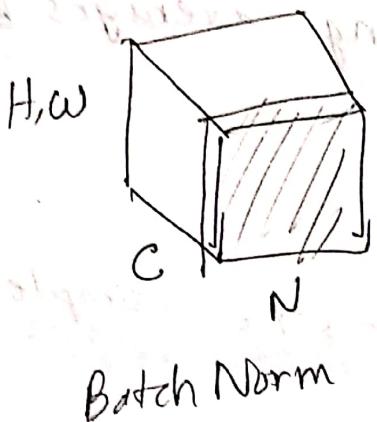
$$y = \frac{\alpha}{\sigma} (x - \mu) + b$$

Learnable bias term

Most norm layers (whiten activations)

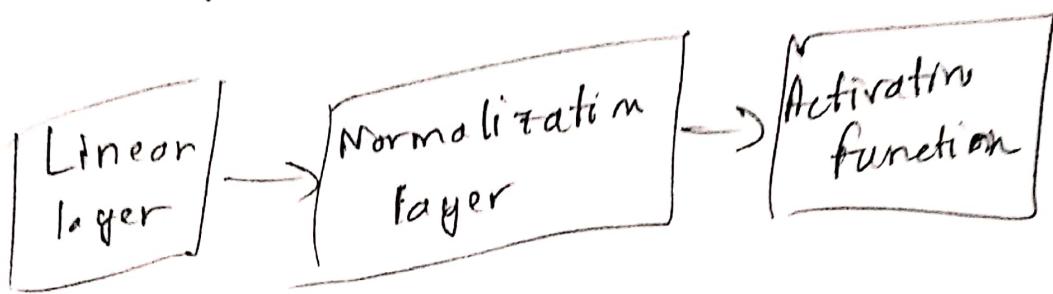
Batch Norm

In practice, group/batch norm work well for computer vision problems and instance/layer norm are heavily used for language problems



- Networks with normalization layers are easier to optimize, allowing for the use of larger learning rates.
- Extra 'noise' results in better generalization.
- Reduce sensitivity to weight initialization.

Group Norm is Recommended since it's more stable
theoretically simple and usually works better



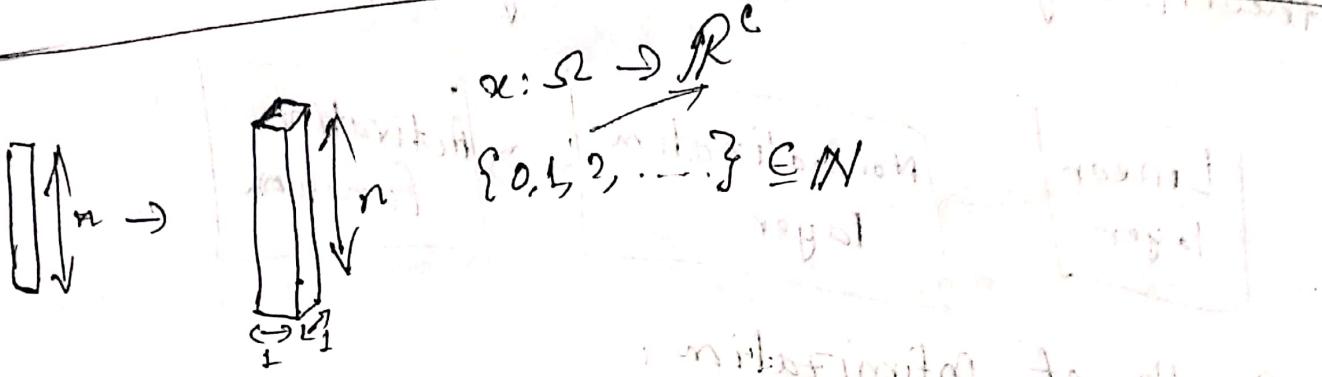
The Death of Optimization:

MRI reconstruction

MRI machines achieve acquiring data one row or column at time of the Fourier domain representation. The inverse Fourier transform is applied at the end to get the anatomical image while processing.

Instead of doing optimization in each step a big neural network can be used to automate this which gives a much more better and fine-grained result.

Practicum 5: 1D multi-channel convolution and autograd



stereophonic signal \rightarrow 2 channels

• mostly done in ~~C/C++~~ code

a. backprop() in pytorch accumulates all gradients backwards w.r.t params

With Great Flexibility Comes Great Responsibility

with torch.no_grad()

[no tracing]

Not possible with tensorflow

Week 6

Lectures

CNN application, RNN and attention

No maximum suppression \rightarrow one highest score is the result

\Rightarrow Word-level training with weak supervision [Matan 1992]

\Rightarrow Multiple Character Recognition [Matan et al. 1992]

Sliding Window ConvNet + Weighted finite state Machine

Face Detection

[Vaillant et al. 93, 94]

ConvNet applied to large images

- ConvNet applied to large images
- Heatmaps at multiple scales
- Non-maximum suppression for candidates
- Non-maximum suppression for 256x256
- 6 second on sparcstation

Semantic Segmentation

[ConvNet for Long Range Adaptive Robot Vision]

[DARPA LAGR Program 2005-2008]

To find out the regions where the robot can drive (Stereo Labeled)

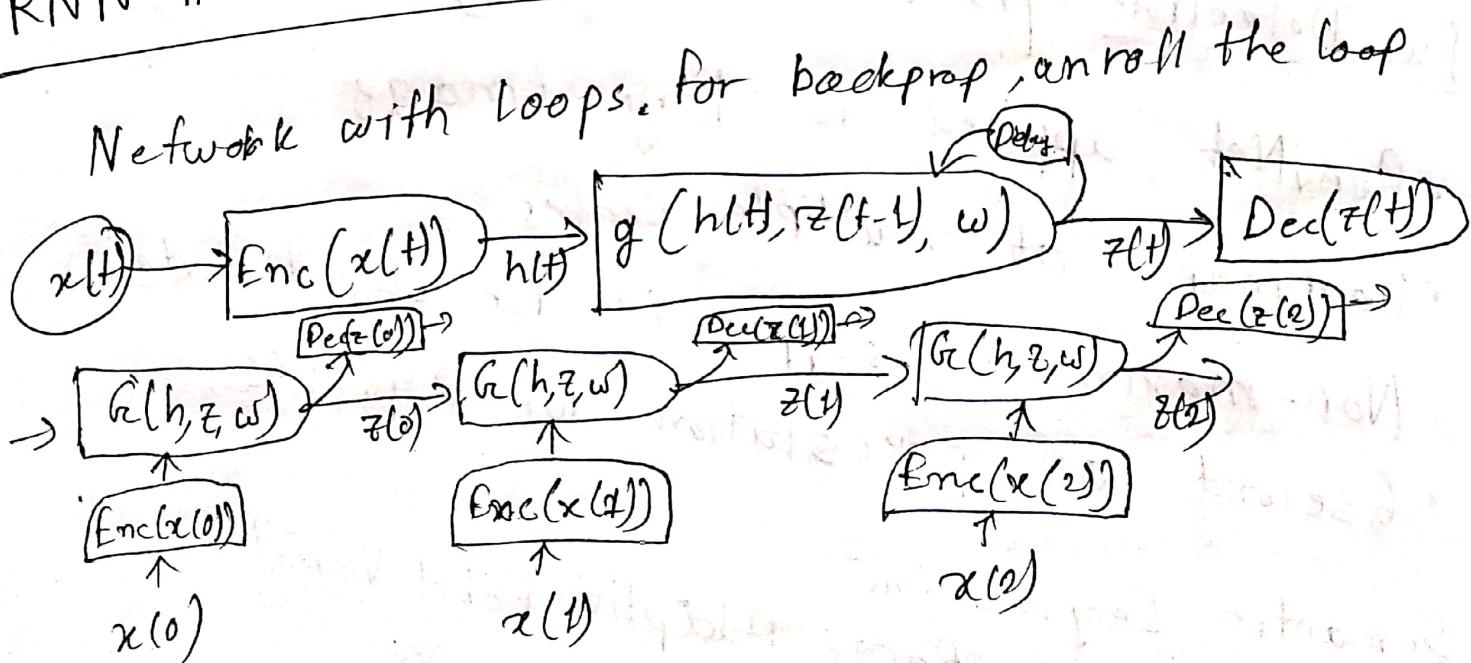
Robot adapted the last layers as it moved on. N-1 layers were frozen and 1 last layer was being trained constantly

lagr@net-scale.com

Scene Parsing / Labeling: Robot fab [Fabre et al
ICML 2013, PAMI 2013]

Multiscale Convnet architecture

RNN Architecture and Attention:



BPTT \rightarrow Backprop Through Time

Vanishing and exploding gradient problem

RNN tricks [Pascanu ICML 2013, Bengio ICASSP 2013]

- Clipping gradients (avoids exploding / gradients)
- Leaky integration (propagate long term dependencies)
- Momentum (cheap 2nd order)
- Initialization (start in right ballpark avoids exploding/vanishing)
- Sparse gradients (symmetry breaking)
- Gradient propagation regularizer (avoid ga vanishing gradient)
- LSTM w/ self-loops (avoid vanishing gradient)

Attention module

$$s_i = \sum_j w_i \cdot x_{ij} \quad w_i = \frac{e^{z_i}}{\sum_k e^{z_k}}$$

GRUs, LSTMs, Seq2Seq

↓
[Hochreiter
&
Schmidhuber 97]

Seq2Seq

Sequence encoder \rightarrow sequence decoder

Applications

\rightarrow Translation [Sutskever NIPS 2014]

\rightarrow Multilayer LSTM

Translation [Bahdanau, Cho, Bengio]

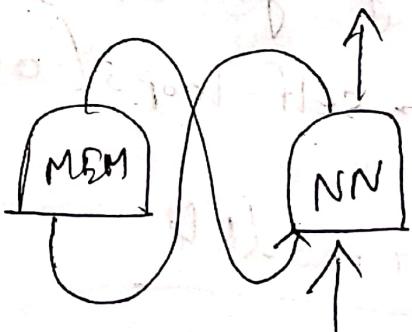
Memory Networks (Story Storing Memory)

Key and values

$$d_i = k_i^T x$$

$$c = \text{softmax}(a)$$

$$s = \sum_i c_i v_i$$



Neural Turing Machine

\rightarrow Differentiable Computation