

Lecture 1 (History, Evolution and revolution of DL) 2/7/2020

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, autopoiesis, 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

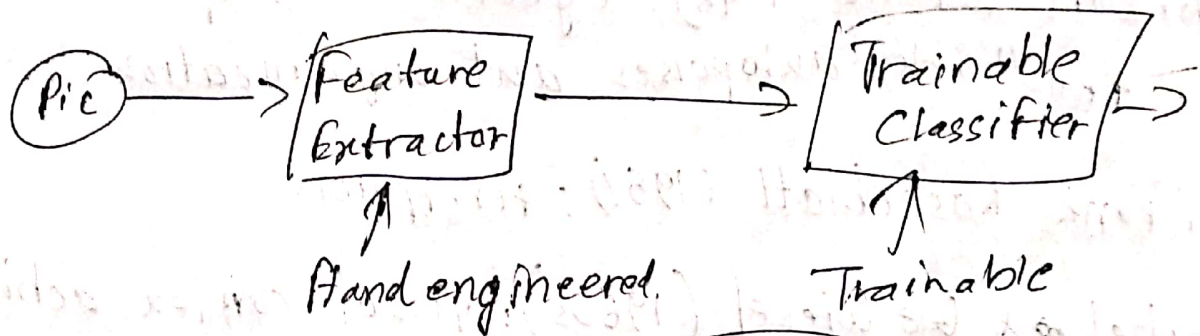
Supervised Learning goes back to the Perceptron and Adaline

The McCulloch-Pitts Binary Neuron

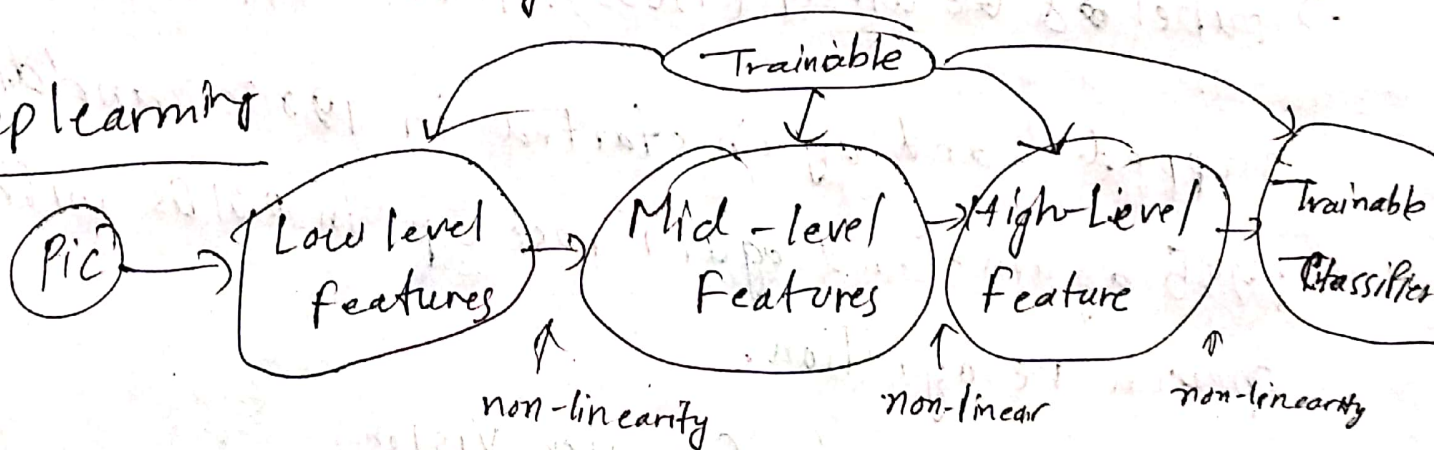
→ Perceptron: Weights are motorized potentiometers

→ Adaline: Weights are electrochemical memristors

Traditional Machine Learning



Deep learning



MLP (Multi-layer Perceptron)

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

Supervised Machine learning = Function Optimization

Stochastic gradient descent

$$W \leftarrow W_i - \eta \frac{\partial L(W, X)}{\partial W_i}$$

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 Cortex.

- simple cells detect local features [1962]

- complex cells "pool" the outputs of simple cells within a 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 [A. Hodsett 2009]

Semantic Segmentation with ConvNets [Farabet 2012]

1986-1996] \Rightarrow ~~the~~ 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]

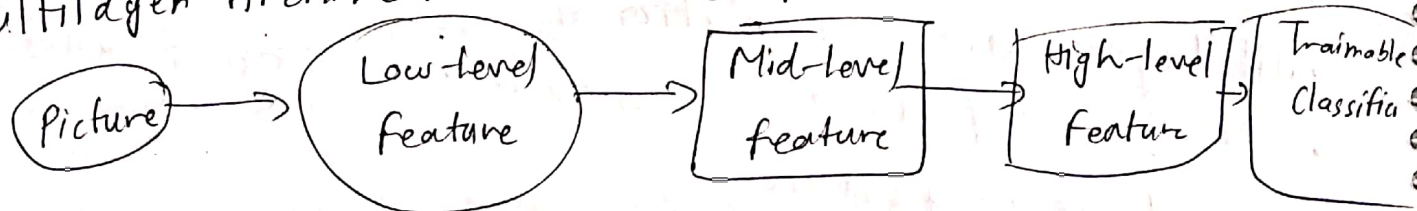
VGG [2013]

GoogLeNet [2013]

ResNet [2015]

DenseNet [2017]

Multilayer Architectures \Rightarrow Compositional structure of Data



MSRA-2015

Mask R-CNN
2017

Mask R-CNN
2017

Mask-RCNN 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

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-linear~~ non-linearity we can't get

Space Stretching

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

Week 2: Lecture: SGD and Backprop

Parameterized Model

$$\hat{y} = G(\underbrace{x}_{\text{params}}, \underbrace{w}_{\text{weights}})$$

[Parameterized Deterministic Function]

Computing function G may involve complicated algo.

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

Deterministic function

Scalar valued function (implicit output)

L 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

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

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

$$y_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

John + Vidal
[vector of probabilities] $\log y_i$

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

Practical tips of Backprop

softmax is just a generalization of
~~set~~ sigmoid with multiple outputs

- Use ReLU (tanh and others are ~~not~~ falling out of favor)
- Use cross entropy for classification
- Use sgd on minibatches
- Shuffle the training samples
- Normalize the input variables (zero mean, unit variance)
- Schedule to decrease the learning rate
- Use a bit of L1 ^{or} ~~and~~ 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

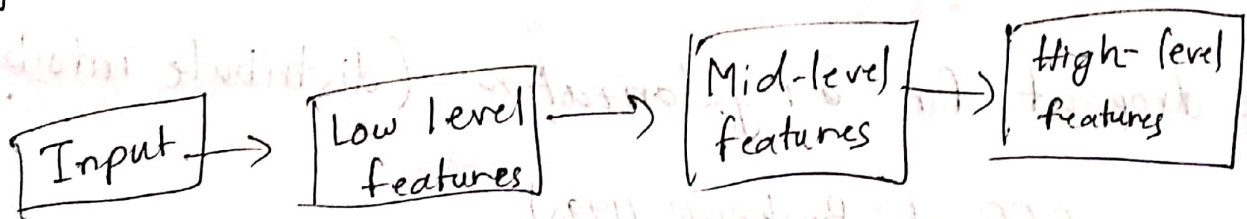
NN are mainly doing rotation and squashing

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 ~~at~~ it is called multi-headed network

Train data

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

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

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

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

1-hot encoding \xrightarrow{k}

$$\psi = \begin{bmatrix} -y^{(1)}- \\ -y^{(2)}- \\ \vdots \\ -y^{(m)}- \end{bmatrix} \quad \begin{matrix} \uparrow \\ m \\ \downarrow \end{matrix}$$

Fully connected layer



$$h = f(w_n x + b_n)$$

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

Affine transformation \rightarrow rotation and scaling

Non-linear function \rightarrow squashing

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

$$b_n \in \mathbb{R}^d$$

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

$$b_y \in \mathbb{R}^k$$

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

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

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

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

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

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

\uparrow cross entropy / negative log likelihood

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

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

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

$$\theta = \{w_n, b_n, w_y, b_y\}$$

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

$$-\frac{d\mathcal{J}(\theta)}{d\theta}(\theta_0)$$

- notebooks
- spiral classification
 - ~~max~~ regression

choice of activation will have effects on regression