Neural Networks: Old and New

Ju Sun

Computer Science & Engineering University of Minnesota, Twin Cities

September 14, 2020

Start from neurons

Shallow to deep neural networks

A brief history of Al

Suggested reading

Model of biological neurons



Credit: Stanford CS231N

Biologically ...

- Each neuron receives signals from its dendrites
- Each neuron outputs signals via its single axon
- The axon branches out and connects via synapese to dendrites of other neurons

Model of biological neurons



A cartoon drawing of a biological neuron (left) and its mathematical model (right).



Mathematically ...

- Each neuron receives x_i 's from its **dendrites**
- x_i 's weighted by w_i 's (synaptic strengths) and summed $\sum_i w_i x_i$
- The neuron fires only when the combined signal is above a certain threshold: $\sum_i w_i x_i + b$
- Fire rate is modeled by an **activation function** f, i.e., outputting $f(\sum_i w_i x_i + b)$

Artificial neural networks



Why called artificial?

- (Over-)simplification on neural level
- (Over-)simplification on connection level

In this course, neural networks are always artificial.

Start from neurons

Shallow to deep neural networks

A brief history of Al

Suggested reading



$$f\left(\sum_{i} w_{i}x_{i} + b\right) = f\left(\boldsymbol{w}^{\mathsf{T}}\boldsymbol{x} + b\right)$$

We shall use σ instead of f henceforth.

Examples of activation function $\boldsymbol{\sigma}$



Credit: [Hughes and Correll, 2016]

One neuron: $\sigma (w^{\intercal}x + b)$ Neural networks (NN): structured organization of artificial neurons



w's and *b*'s are unknown and need to be learned Many models in machine learning **are** neural networks

Supervised Learning

- Gather training data $(oldsymbol{x}_1,oldsymbol{y}_1),\ldots,(oldsymbol{x}_n,oldsymbol{y}_n)$
- Choose a family of functions, e.g., \mathcal{H} , so that there is $f \in \mathcal{H}$ to ensure $\boldsymbol{y}_i \approx f\left(\boldsymbol{x}_i\right)$ for all i
- Set up a loss function ℓ to measure the approximation quality
- Find an $f \in \mathcal{H}$ to minimize the average loss

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} \ell\left(\boldsymbol{y}_{i}, f\left(\boldsymbol{x}_{i}\right)\right)$$

 \ldots known as **empirical risk minimization** (ERM) framework in learning theory

A typical setup

Supervised Learning from NN viewpoint

- Gather training data $(oldsymbol{x}_1,oldsymbol{y}_1)\,,\ldots,(oldsymbol{x}_n,oldsymbol{y}_n)$
- Choose a NN with k neurons, so that there is a group of weights, e.g., $(\boldsymbol{w}_1,\ldots,\boldsymbol{w}_k,b_1,\ldots,b_k)$, to ensure

$$\boldsymbol{y}_i \approx \{ \mathsf{NN} (\boldsymbol{w}_1, \dots, \boldsymbol{w}_k, b_1, \dots, b_k) \} (\boldsymbol{x}_i) \quad \forall i$$

- Set up a loss function ℓ to measure the approximation quality
- Find weights $(\boldsymbol{w}_1,\ldots,\boldsymbol{w}_k,b_1,\ldots,b_k)$ to minimize the average loss

$$\min_{\boldsymbol{w}'s, b's} \frac{1}{n} \sum_{i=1}^{n} \ell \left[\boldsymbol{y}_i, \left\{ \mathsf{NN}\left(\boldsymbol{w}_1, \dots, \boldsymbol{w}_k, b_1, \dots, b_k \right) \right\} (\boldsymbol{x}_i) \right]$$

Linear regression



– Data: $(oldsymbol{x}_1, y_1), \ldots, (oldsymbol{x}_n, y_n)$, $oldsymbol{x}_i \in \mathbb{R}^d$

- Model:
$$y_i \approx w^\intercal x_i + b$$

- Loss: $||y \hat{y}||_2^2$
- Optimization:

$$\min_{\boldsymbol{w},b} \frac{1}{n} \sum_{i=1}^{n} \|y_i - (\boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_i + b)\|_2^2$$

Credit: D2L



 $\boldsymbol{\sigma}$ is the identity function

Perceptron



Frank Rosenblatt

(1928-1971)

- Data: $(\boldsymbol{x}_1, y_1), \dots, (\boldsymbol{x}_n, y_n)$, $\boldsymbol{x}_i \in \mathbb{R}^d$, $y_i \in \{+1, -1\}$
- Model: $y_i \approx \sigma \left(\boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_i + b \right)$, σ sign function



- Loss: $\mathbf{1}\left\{y\neq\hat{y}
 ight\}$
- Optimization:

$$\min_{\boldsymbol{w}, b} \frac{1}{n} \sum_{i=1}^{n} \mathbf{1} \left\{ y_i \neq \sigma \left(\boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_i + b \right) \right\}$$

Perceptron

Perceptron is a single artificial neuron for binary classification



dominated early AI (50's - 70's)

Logistic regression is similar but with sigmod activiation



Softmax regression

- Data: $(x_1, y_1), \ldots, (x_n, y_n)$, $x_i \in \mathbb{R}^d$, $y_i \in \{L_1, \ldots, L_p\}$, i.e., multiclass classification problem
- Data preprocessing: labels into vectors via one-hot encoding

$$L_k \Longrightarrow [\underbrace{0, \dots, 0}_{k-1 \, 0's}, 1, \underbrace{0, \dots, 0}_{n-k \, 0's}]^{\mathsf{T}}$$

So: $y_i \Longrightarrow \boldsymbol{y}_i$

- Model: $y_i \approx \sigma (W^{\mathsf{T}} x_i + b)$, here σ is the softmax function (maps vectors to vectors): for $z \in \mathbb{R}^p$,

$$\boldsymbol{z} \mapsto \left[\frac{e^{z_1}}{\sum_j e^{z_j}}, \dots, \frac{e^{z_p}}{\sum_j e^{z_j}}
ight]^\mathsf{T}$$

- Loss: cross-entropy loss $-\sum_j y_j \log \hat{y}_j$
- Optimization ...

... for multiclass classification



Multilayer perceptrons



Also called feedforward networks or fully-connected networks

Modern NNs: many hidden layers (deep), refined connection structure and/or activations

- Linear regression
- Perception and Logistic regression
- Softmax regression
- Multilayer perceptron (feedforward NNs)
- Support vector machines (SVM)
- PCA (autoencoder)
- Matrix factorization

see, e.g., Chapter 2 of [Aggarwal, 2018].

Start from neurons

Shallow to deep neural networks

A brief history of AI

Suggested reading

Birth of AI



- Crucial precursors: first computer, Turing test
- 1956: Dartmouth Artificial Intelligence Summer Research
 Project Birth of AI



Turing Test



Alan Turing (1912-1954)

First golden age



Symbolic AI: based on rules and logic





rules for recognizing dogs?

First AI winter



Gartner hype cycle

Perceptron



invented 1962



written in 1969, end of Perceptron era



Marvin Minsky (1927–2016)

Birth of computer vision

MASSACHUSETTS INSTITUTE OF TECHNOLOGY PROJECT MAC

Artificial Intelligence Group Vision Memo. No. 100. July 7, 1966

THE SUMMER VISION PROJECT

Seymour Papert

The summer vision project is an attempt to use our summer vorkers effectively in the construction of a significant part of a visual system. The particular task was chosen partly because it can be segmented into sub-problems which will allow individuals to work independently and yet participate in the construction of a system complex enough to be a real lumbark in the development of "mattern recognized."

1966

VISION

Convrighted Material



David Marr

FOREWORD BY Shimon Ullman AFTERWORD BY Tomaso Poggio

Copyrighted Material

around 1980

Second golden age



expert system





Can we build comprehensive knowledge bases and know all rules?

Key ingredients of DL have been in place for 25-30 years:

Landmark	Emblem	Epoch
Neocognitron	Fukushima	1980
CNN	Le Cun	mid 1980s'
Backprop	Hinton	mid 1980's
SGD	Le Cun, Bengio etc	mid 1990's
Various	Schmidhuber	mid 1980's
CTF	DARPA etc	mid 1980's

After 2nd Al winter



Machine learning takes over ...



Starting 1990's

Support vector machines (SVM) Adaboost Decision trees and random forests Deep learning Start from neurons

Shallow to deep neural networks

A brief history of Al

Suggested reading

- Chap 2, Neural Networks and Deep Learning.
- Chap 3-4, Dive into Deep Learning.
- Chap 1, Deep Learning with Python.

- [Aggarwal, 2018] Aggarwal, C. C. (2018). Neural Networks and Deep Learning. Springer International Publishing.
- [Hughes and Correll, 2016] Hughes, D. and Correll, N. (2016). Distributed machine learning in materials that couple sensing, actuation, computation and communication. *arXiv:1606.03508*.