

Artificial Intelligence in data science

Introduction

Janos Török

Department of Theoretical Physics

September 8, 2022

Information

► Coordinates:

- Török János
- Email: torok.janos@ttk.bme.hu, torok72@gmail.com
- Consultation:
 - F III building, first floor 6 (after the first stairs to the right, at the end of the corridor), Department of Theoretical Physics
 - Upon demand (Email)

► Webpage:

https://physics.bme.hu/BMETE15MF75_kov?language=en

► Homework, project: Moodle at <https://edu.ttk.bme.hu/>

► Login to moodle: University login

► Teams: Two years ago I have recorded the lecture parts which are available in the teams. Please, do not use the teams messaging system to communicate with me. Use email! I send you teams code upon request.

Required knowledge

- ▶ Knowledge of basic matrix and linear algebra
- ▶ Python language, numpy, matrices, data handling
- ▶ **Please note:** If today's class poses serious difficulties to you then I suggest dropping this course

Used software

- ▶ Python3
- ▶ Suggested tools
 - ▶ google colab: Go to <https://colab.research.google.com>, or save an .ipynb notebook in your google drive and open it. google will offer to open it with colab
 - ▶ Anaconda: <https://docs.anaconda.com/>, be careful, installation of some packages, e.g. tensorflow may require some expertise. Also note that I do not use Windows at all, so I will not be able to help with the installation, but the Internet may!
- ▶ Book: Stuart J. Russell and Peter Norvig: *Artificial Intelligence A Modern Approach*, uploaded in the teams group in Files under the General channel

Requirements

▶ Pass

- ▶ 50% homework submitted and accepted
- ▶ Project presented and accepted

▶ Mark

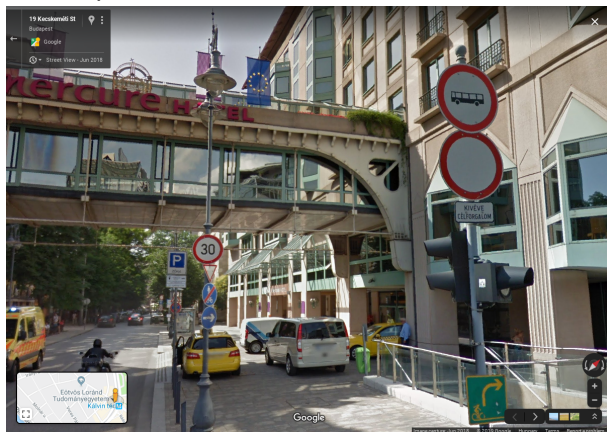
- ▶ **100 points/HW**: Homeworks (individual)
- ▶ **500 points**: Data challenge (pairs)
- ▶ **1000 points**: Project (pairs) presented at the end of the semester
- ▶ Mark:
 1. –1099
 2. 1100 – 1399
 3. 1400 – 1699
 4. 1700 – 1999
 5. 2000–
- ▶ At the end of the semester there will be a 5 minute investigation to verify the authenticity of the codes
- ▶ Turn it in language: English, Hungarian, German, French

Submission of works

- ▶ Using fancy visualization techniques does not improve the mark
- ▶ Use markdown to explain important steps in few words.
- ▶ Submit always the notebook with the outputs, not the pdf version of it, however you may submit an extra pdf if you want to give more explanation. (zip them into a single file.)
- ▶ Please keep the outputs in the notebooks before submitting. If you still run into size problems in the moodle (it happened with some of the final projects, then remove the big outputs, but keep the ones that may help me evaluate you works.
- ▶ The reason behind this request is that there are 28 students, 12 homeworks, plus a project, so I am supposed to correct ~ 200 works, which is a lot.

Data challenge: possibility 1

- ▶ Traffic sign recognition
- ▶ Images from street view
- ▶ Only circular red ones, into four categories.
- ▶ Test random street from Budapest, score based on accuracy.
- ▶ Expected success rate $\sim 30\%$



Data challenge: possibility 2

- ▶ Covid data prediction, in light of government measures
- ▶ Task will require data collection and cleaning
- ▶ Task will be the prediction the effectiveness of the counter measures

Topics

- ▶ Image segmentation
- ▶ Decision tree
- ▶ Deep learning (from scratch in numpy)
- ▶ Higher level implementations (tensorflow, sklearn)
- ▶ Sequential data
- ▶ Reinforcement learning, Game models
- ▶ Textual data
- ▶ Convolutional neural networks
- ▶ Pre trained models, merging models
- ▶ Data augmentation, autoencoding

Aim

- ▶ Introduction to machine learning
- ▶ Physicist's view: learn how it works and leave the tricks for engineers
- ▶ Simple examples of different machine learning methods, while trying to avoid **big magic happens here!** moments
- ▶ You will learn how to apply different type of neural networks, different learning methods, and how they work
- ▶ What I cannot teach you (I am not an expert): What are the best parameters for certain cases, how to optimize learning speed, efficiency, etc.
- ▶ No big scale beautiful examples :-)

Artificial Intelligence

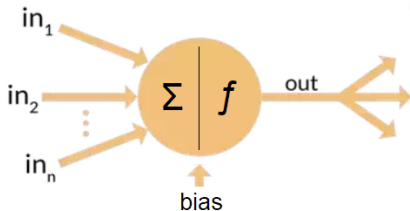
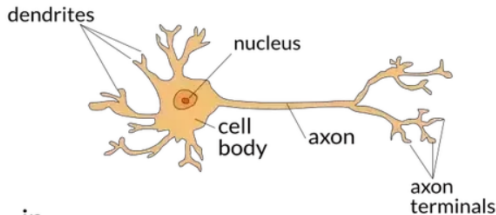
- ▶ "a system's ability to correctly interpret external data, to learn from such data, and to use those learnings to achieve specific goals and tasks through flexible adaptation."
- ▶ "AI is whatever hasn't been done yet."
- ▶ Weak or strong
 - ▶ Weak: trained for one particular task (e.g. Google assistant)
 - ▶ Strong: Can perform unfamiliar task.

Artificial Intelligence types

- ▶ Reactive machines: single purpose machines trained on data, e.g. Deep blue
- ▶ Limited memory: uses past (limited extent) to make decisions in the present
- ▶ Theory of mind: understand others and their decisions (not yet)
- ▶ Self-awareness: What everybody is afraid of!

Artificial neuron

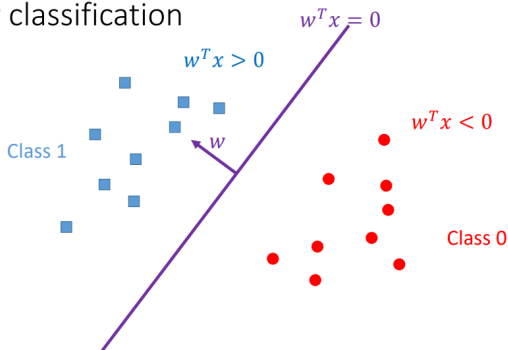
- ▶ Biological neuron:
 - ▶ Stimulus in dendrites
 - ▶ Fire (activate axon) when stimulus is large enough
- ▶ Artificial neuron, perceptron:
 - ▶ Weighted sum of the input
 - ▶ Output is a non-linear function of the input



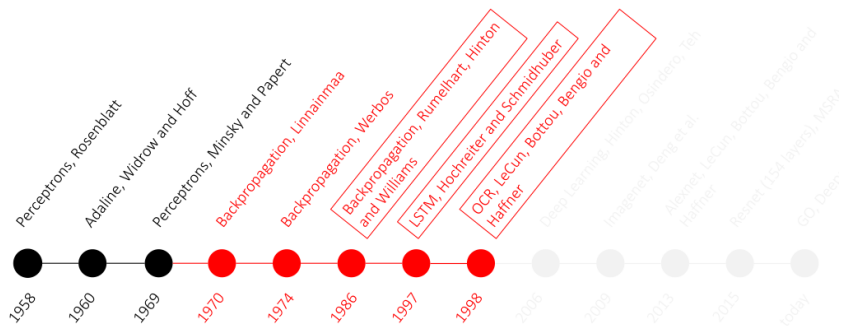
Classification by a perceptron

- ▶ Similar to linear or logistic regression
- ▶ Cuts the input space along a line

Linear classification



AI winter: spite of discoveries



Problems:

- ▶ Cannot solve the XOR problem, why bother
- ▶ Far fetched expectations (best sci fi era)
- ▶ Insufficient computational power
- ▶ Cuts in funding
- ▶ However significant new discoveries

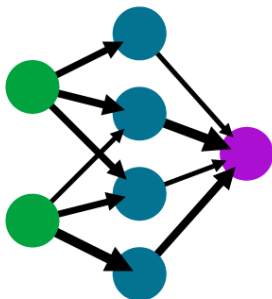
Source: UVA deep learning course - Efstratios Gavves

Deep neural networks

- ▶ Cut the space by more nonlinear lines
- ▶ Stack the neurons
 - ▶ Input vector I
 - ▶ Hidden layers
 - ▶ Output vector $O(I)$
 - ▶ Transition matrix W_{ij}
 - ▶ Learning using a cost function

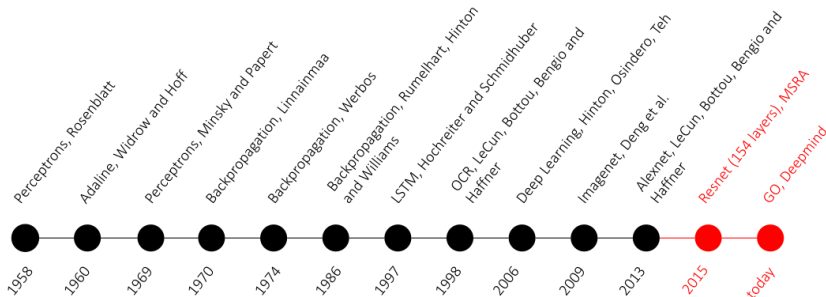
A simple neural network

input layer hidden layer output layer



Deep neural networks: the breakthrough

- ▶ Learning methods: layer-by-layer
- ▶ Learning algorithms: SGD, Adam
- ▶ Datasets: 2009 Imagenet dataset 16 million images annotated by humans in hierarchical categories
- ▶ Computational power, especially GPU
- ▶ Theory improvements: ReLu, dropout, data augmentation, regularization



Neural networks: Learning

- ▶ Supervised learning
- ▶ Data training:
 - ▶ Supervised learning
 - ▶ Fitness function, energy:

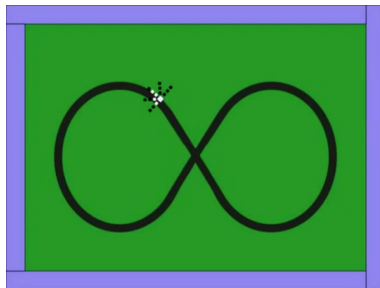
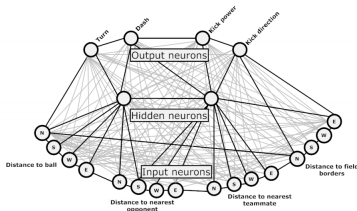
$$E = |T(I) - O(I)|,$$

where $T(I)$ is the target vector for input I

- ▶ Minimize E for available set of $\{I, I(O)\}$ pairs
 - ▶ Deep learning: many layers of neurons in the neural network
- ▶ Test goodness:
 - ▶ Use only part of $\{I, I(O)\}$ pairs for learning, the rest is for testing.
- ▶ Used for: pattern recognition, classification, etc.

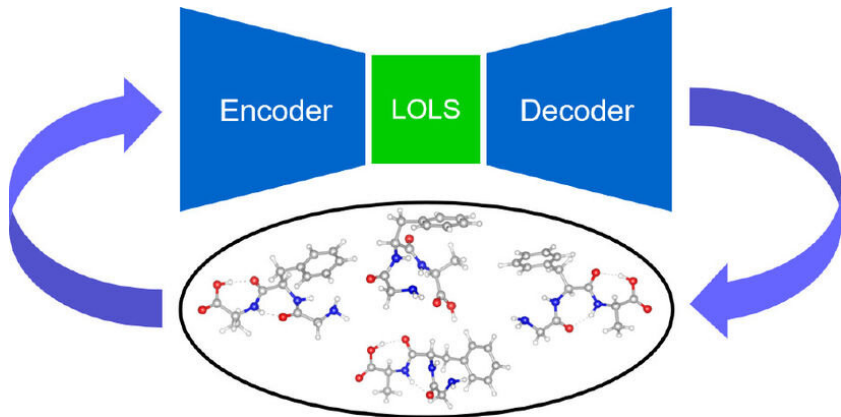
Neural networks: Learning

- ▶ Reinforcement learning
- ▶ Cost function is a long time performance on an agent making decisions based on the neural network.
- ▶ Test goodness:
 - ▶ Compare with other agents which can be algorithmic or based on neural networks
- ▶ Used for: control problems, AI, complex optimization



Neural networks: Learning

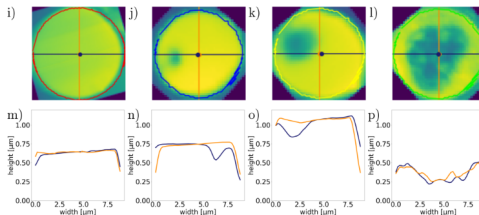
- ▶ Unsupervised learning, weight is increased for neurons that fire together
- ▶ Autoencoding, etc.



Neural networks: Models

Which model to use?

- ▶ Human created models
 - ▶ No need for large amount of data
 - ▶ Less memory and CPU requirements (in most cases)
 - ▶ Example: red blood cell with malaria, asymmetric cut



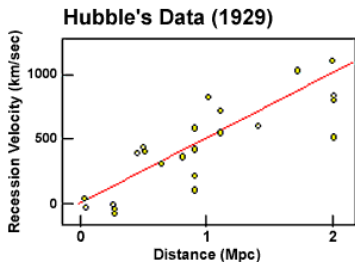
- ▶ Neural network model?
 - ▶ Let the algorithm find the features
 - ▶ Human effort is to be put in better AI algorithm and not in use case specific applications
 - ▶ Huge memory and C/GPU requirements

Linear regression

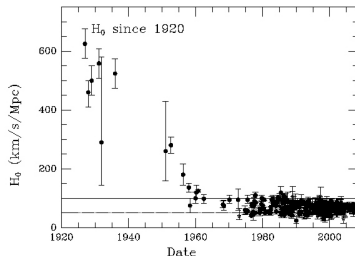
- ▶ y is a linear function of the parameters x
- ▶ Linear fit in m dimensions
- ▶ For n data points
- ▶ Surprisingly efficient

Linear regression

Hubble constant



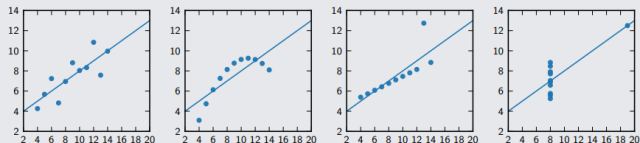
change in time:



Linear regression

EXAMPLE: ANSCOMBE'S QUARTET REVISITED

Recall Anscombe's Quartet: 4 datasets with very similar statistical properties under a simple quantitative analysis, but that look very different. Here they are again, but this time with linear regression lines fitted to each one:



For all 4 of them, the slope of the regression line is 0.500 (to three decimal places) and the intercept is 3.00 (to two decimal places). This just goes to show: visualizing data can often reveal patterns that are hidden by pure numeric analysis!

Solution: e.g. random sampling of points, and stability of solution

MIT: unknown lecturer

Linear regression

- ▶ Assume linear form for the loss function:

$$f_w(x) = w^T x$$

- ▶ i indexes data points, m indexes parameters
- ▶ The problem must not be linear, e.g. polynomial fit w contains the coefficients of the polynomial:

$$y_i = w_0 + w_1 x_i + w_2 x_i^2 + \dots + w_m x_i^m + \varepsilon_i$$
$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & x_1 & x_1^2 & \dots & x_1^m \\ 1 & x_2 & x_2^2 & \dots & x_2^m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^m \end{pmatrix} \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

- ▶ If X is a square matrix than $w = X^{-1}y$
- ▶ Otherwise $w = (X^T X)^{-1} X^T y$

Linear regression

- ▶ If X is a square matrix than $w = X^{-1}y$
- ▶ Otherwise $w = (X^T X)^{-1} X^T y$

