Deep learning, dream or reality?

Roberto Battiti University of Trento (Italy) Lobachevsky State University of Nizhni Novgorod



SCIENTIFIC YOUTH SCHOOL HIGH-PERFORMANCE COMPUTING, OPTIMIZATION AND APPLICATIONS

Lobachevsky State University of Nizhni Novgorod

Learning from Data and Machine Learning



If you show a picture to a three-year-old and ask if there is a tree in it, you will likely get the correct answer. If you ask a thirty-year-old what the definition of a tree is, you will likely get an inconclusive answer. We didn't learn what a tree is by studying the mathematical definition of trees. We learned it by looking at trees. In other words, we learned from 'data'.

Three ways of model building

output



Model

Useful for what-if questions predictions Input optimization

input

1) Explicit and rigid models

Pressure = N k T / V



(Volume, Temperature)

e.g., Physics: Boyles's law:

"For a fixed mass of gas, at a constant temperature, the product (pressure x volume) is a constant."

PV = N k T

2) Parametric, with statistics

Output voltage







Ronald Fisher in 1913



e.g., Maximum likelihood estimation

Input voltage

3) Non-parametric models, neural nets, modern ML (1960++, 1985, 2010)



(Movie, Viewer)

Eduardo Caianiello, 1961

Ivakhnenko, Alexey (1965). Cybernetic Predicting Devices. Kiev: Naukova Dumka.



Very flexible, no rules elicitation, Only need abundant (relevant) data

The dream

"give computers the ability to learn without being explicitly programmed" (<u>Arthur Samuel</u>, 1959).

The Tool

Weights of the flexible model are determined via optimization, but aiming at generalization (learning is *mean* not *end*)

Movies and Viewers

- Movie1 = (1.2, 3.3, 2.1, ..., ..., 7.7)
- Movie2 = (3.2, 5.6, 1.2, ..., ..., 3.4)
- Viewer1 = (6.2, 5.6, 7.2, 2.1)

Map to vectors of the same dimensions \rightarrow m, v

Obtain recommandation by simple scalar product

Objective = Sum_data_i (**m_i . v_i -** r_i)²

Minimize to determine vectors!

Movies and Viewers



Setup of the Machine Learning problem





Linear models

Most right-handed people are linear thinking, think inside the box.



Linear models

• Just below the mighty power of optimization lies the awesome power of linear algebra.



Data about price and power of different car models. A linear model (fit) is shown.

Linear regression

A linear dependence of the output from the input features

```
f(x) = w_1 x_1 + w_2 x_2 + ... + w_d x_d.
```

- The model is simple, can be easily trained,
- The computed weights in the linear summation provide a direct explanation of the importance of the various attributes

Best (linear) fit, by Least Squares → optimization

 Errors can be present (every physical quantity can be measured only with a finite precision)

 $y_i = \boldsymbol{w}^T \cdot \boldsymbol{x}_i + \boldsymbol{\epsilon}_i,$

- Determine optimal weight vector \mathbf{w} so that: $\hat{f}(\mathbf{x}) = \mathbf{w}^T \cdot \mathbf{x}$ approximates as closely as possible the experimental data
- minimizes the sum of the squared errors (least squares approximation):

ModelError
$$(w) = \sum_{i=1}^{c} (w^T \cdot x_i - y_i)^2.$$

Biological motivations of linear models



Neurons and synapses in the human brain

Abstract model: the perceptron



• Output is a weighted sum of the inputs passed through a final threshold function.

Minimizing the sum of squared errors

- If zero measurement errors and a perfect linear model:
 set of linear equations w^T x_i = y_i one for each example
- in real-world cases, reaching zero for the ModelError is impossible, and the number of data points can be much larger than the number of parameters d.
- furthermore, the goal of learning is *generalization* (and requiring zero error can cause "overtraining")

Minimizing the sum of squared errors

• Error is quadratic in parameters **w**

ModelError
$$(w) = \sum_{i=1}^{\ell} (w^T \cdot x_i - y_i)^2.$$

- Find minimum:
 - take partial derivatives
 - equate them to 0
- Obtain linear equations, typically more equations than examples

Minimizing the sum of squared errors

From inverse to pseudo-inverse, solution is:

$$\boldsymbol{w}^* = (X^T X)^{-1} X^T \boldsymbol{y};$$

where $y = (y_1; ...; y_L)$ and **X** is the matrix whose rows are the **x**_i vectors.

- least-square and pseudo-inverse are among the most popular tools in Science and Engineering
- alternative is gradient descent

An analogy in Physics



Spring analogy for least squares fits.

Springs connect a rigid bar to the experimental points.

The best fit is the line that **minimizes the overall potential energy of the system** (proportional to the sum of the squares of the spring length).

Neural networks with multiple layers

Quegli che pigliavano per altore altro che la natura, maestra de' maestri, s'affaticavano invano. (Leonardo Da Vinci)



Counter-example: XOR function

Cannot separate points with a line in two dim.



But points *can* be separated by a plane after mapping them into 3D

The biological metaphor

- Our neural system is composed of 100 billion computing units (neurons) and 10¹⁵ connections (synapses).
- How can a system composed of many simple interconnected units give rise to highly complex activities?
- Emergence: complex systems arise out of a multiplicity of relatively simple *interacting* units.



Drawings of cortical lamination by Santiago Ramon y Cajal, each showing a vertical cross-section, with the surface of the **cortex** at the top. The different stains show the cell bodies of neurons and the dendrites and axons of a random subset of neurons.

Symbolic vs. sub-symbolic paradigm

- "Standard" sequential computers operate in cycles, fetching items from memory, applying mathematical operations and writing results back to memory.
- The *intelligence* of biological brains is different, it lies in the interconnection strengths, learning occurs by modifying connections (dynamical systems)
- Neural networks do not separate memory and processing but operate via the flow of signals through the network connections.

Artificial Neural Networks

- A neuron is modeled as a simple computing unit, a scalar product w x ("pattern matching") followed by a sigmoidal ("logistic") function.
- The complexity comes from having more interconnected layers of neurons involved in a complex action (if linear layers are cascaded, the system remains linear)
- The "squashing" functions is essential to introduce nonlinearities in the system

Multilayer Perceptrons (MLP)

- By applying a sigmoidal transfer function to the unlimited output of a linear model, one obtains an output in [0,1] which can be interpreted as a probability
- For classification, hyperplane boundaries become *"fuzzy"*





Multilayer Perceptrons (MLP)(2)

- Composing linear transformations → still linear functions.
- Composing linear functions with nonlinear sigmoids one can approximate all smooth functions. One hidden layer is sufficient.
- The first linear transformation provides a first "hidden layer" of outputs, the second transformation produces the visible outputs from the hidden layer.

MLP architecture

 MLPs are composed of a large number of interconnected units working in parallel and organized in layers with a feedforward information flow.



MLP architecture (2)

- The signals flow sequentially from the input to the output layer.
- For each layer, each unit does the following:
- scalar product between a vector of weights and the vector of outputs of the previous layer;
- 2. nonlinear function to each result to produce the input for the next layer
- A popular transfer function is the sigmoidal (or "logistic") function



MLPs are universal approximators

- What is the flexibility of the MLP architecture to represent input-output mappings?
- An MLP with one hidden layer and a sufficient number of hidden nodes can approximate any smooth function to any desired accuracy.
- MLPs are very flexible ("non-parametric") models: they can approximate any smooth inputoutput transformation.



Analyzing a neural network output with LIONoso Sweeper. The output value, the energy consumed to heat a house in winter, is shown as a function of input parameters. Color coded output (left), surface plot (right). Nonlinearities are visible.

Error Backpropagation

- How do we **learn** optimal MLPs from examples?
- 1. take a "guiding" function to be optimized (e.g., sumof-squared errors on the training examples)
- 1. Use gradient descent with respect to the weights to find the better and better weights
- Stop the descent when results on a validation set are best (if over-learning, generalization can worsen). Learning is not an end, but a *means* for generalizing.

Backpropagation(2)

- One needs derivatives, a simple analysis exercise.
- MLP is a composition of squash functions and scalar products.
- Derivatives can be calculated by using the chain rule for derivatives of composite functions.
- Complexity is O(*number of weights*).
- Formulas are similar to those used for the forward pass, but going in contrary direction, hence the term error backpropagation.
- After the network is trained, calculating the output from the inputs requires a number of simple operations proportional to the number of weights.

Batch backpropagation

• Given an MLP, define its sum-of-squareddifferences energy as:

$$E(w) = \frac{1}{2} \sum_{p=1}^{P} E_p = \frac{1}{2} \sum_{p=1}^{P} (t_p - o_p(w))^2$$

- 1. Let the initial weights be randomly distributed Partial derivatives
- 2. Calculate the gradient $g_k = \nabla \tilde{E}(w_k)$
- 3. The weights at the next iteration k + 1 are updated as follows $w_{k+1} = w_k \epsilon g_k$.

Bold-Driver Backpropagation

- How do we select the learning rate ε ? If small, the learning time increases, if big, the energy oscillates wildly.
- If successive steps reduce E(w), ε increases exponentially
- 2. If E(w) increases, ε decreases rapidly

$$\epsilon(t) = \begin{cases} \rho \ \epsilon(t-1), & \text{if } E(w(t)) < E(w(t-1)); \\ \sigma^l \ \epsilon(t-1), & \text{if } E(w(t)) > E(w(t-1)) \end{cases}$$

 ρ is close to one (1.1) in order to avoid frequent "accidents" σ is chosen to provide a rapid reduction(0.5), and I is the minimum integer such that the reduced rate succeeds in diminishing the energy

On-line or stochastic backpropagation

- E is a sum of many terms, one for each pattern p
- The gradient is a sum of the corresponding partial gradients $\nabla E_p(w_k)$
- In "batch" gradient descent, *first* the contributions $\nabla E_p(w_k)$ are summed, *then* the small step is taken
- Stochastic BP: randomly choose a pattern p and take a small step along a single negative $\nabla E_p(w_k)$ *immediately* after calculating it.

On-line or stochastic backpropagation (2)

• Stochastic on-line backpropagation update:

 $w_{k+1} = w_k - \epsilon \nabla E_p(w_k),$

- where the pattern p is chosen randomly from the training set at each iteration and ε is the learning rate
- Pros: using partial gradients is faster
- Cons: less guarantee of convergence

Advanced optimization for MLP training

- Higher-order derivatives can be used to enhance the search.
- Conjugate gradient and secant methods update an approximation of the Hessian by using only gradient information.
- They are useful for problems with few weights (approx < 100) and requiring high precision in the output value.

Statistical Learning Theory and Support Vector Machines (SVM)



prof. Vapnik

Linear Models: safe separation with margin maximization



SVN for Linearly Separable Problems

• Margin maximization:

 $\begin{array}{ll} \text{Minimize}_{\boldsymbol{w},b} & \frac{1}{2} \|\boldsymbol{w}\|^2\\ \text{subject to} & y_i(\boldsymbol{w}\cdot\boldsymbol{x}_i+b) \geq 1 \quad i=1,\ldots,\ell. \end{array}$

 After introducing Langrange multipliers for the constraints: Quadratic Programming

> Maximize_A $\mathbf{\Lambda} \cdot \mathbf{1} - \frac{1}{2}\mathbf{\Lambda} \cdot D \cdot \mathbf{\Lambda}$ subject to $\begin{cases} \mathbf{\Lambda} \cdot \mathbf{y} = 0 & ;\\ \mathbf{\Lambda} \ge 0 \end{cases}$

SVM Intuition



Figure 12.3: Hypothesis space constraint. The separating hyperplane must maximize the margin. Intuitively, no point has to be too close to the boundary so that some noise in the input data and future data generated by the same probability distribution will not ruin the classification.

Nonlinear Problems: Linear Combination of Kernel functions

$$\sum_{i=1}^{\ell} y_i \lambda_i^* K(x, x_i),$$
$$K(x, x_i) = \varphi(x) \cdot \varphi(x_i).$$

No need to compute

Linear Combination of Kernel functions



Easily solved with **Quadratic Programming**

But need to determine proper Kernel (by hand !)

SVM

- A long period of success
- Statistical Learning Theory (conditions for learning)
- But practical difficulties to determine the Kernel (a lot of work is done by *expert people*, the dream of learning machines is not fully realized).

Deep neural networks

- Some classes of input-output mappings are easier to build if more hidden layers are considered.
- The dream: feed examples to an MLP with many hidden layers and have the MLP automatically develop internal representations (encoded in the activation patterns of the hidden-layers).

Deep Learning



- Bengio, Yoshua (2009). <u>"Learning Deep Architectures for AI"</u> (PDF). Foundations and Trends in Machine Learning. 2 (1): 1– 127.
- Hinton, G. E.; Srivastava, N.; Krizhevsky, A.; Sutskever, I.; Salakhutdinov, R.R. (2012). "Improving neural networks by preventing co-adaptation of feature detectors«
- Bengio, Yoshua; LeCun, Yann; Hinton, Geoffrey (2015). "Deep Learning". Nature. 521: 436–444.

Practical obstacles to deep MLPs

- Partial derivatives w.r.t. the weights of the first layers tend to be very small, leading to numerical estimation problems.
- As a result, it can happen that the internal representations developed by the first layers will not differ too much from being randomly generated, and leaving only the topmost levels to do some "useful" work.
- A very large number of parameters (such as in deep MLP) can lead to overtraining

Deep learning

Lately, deep learning has lead to superior classification results in challenging areas

The main scheme is as follows:

1. use unsupervised learning from many unlabeled examples to prepare the deep network in an initial state;

2. use back propagation only for the final tuning with the set of labeled examples

+ use tricks like convolutional networks (sharing of parameters) and proper architecture to encourage invariance

Deep Learning



Feature detectors in a frog retina (*Bufo Bufo*) are hard-wired and specialized to detect a fly at the distance that the frog could

Convolutional Neural Networks



Convolutional Neural Networks



Auto-encoders

- Auto-encoders build internal representations in an unsupervised manner
- One builds a network with a hidden layer and demands that the output simply reproduces the input
- Squeezing: in the hidden layer the input is compressed into an encoding c(x) with less variables than the original one
- c(x) will be forced to discover regularities in the input patterns

Auto-encoders (2)



Auto-encoders (2.1)



The codes produced by a 2000- 500-250-125-2 autoencoder on news stories by Reuters. Clusters corresponding to different topics, with different colors, are clearly visible (details in [187]).

Auto-encoders (3)

- The auto-encoder can be trained by backpropagation
- Classification labels are not necessary
- After the auto-encoder is built, the hidden layer (weights and hidden units) is transplanted to a second network with an additional layer, intended for classification
- This new network will be trained on the **labelled** examples to realize a classifier.

Auto-encoders (4)



Auto-encoders (5)

 A chain of subsequent hidden layers by iterating and composing subsequent encodings c'(c(x))

At each iteration, the auto-encoder derives a more compressed internal representation



 Appropriate numbers for the number of layers and the optimal number of units can be obtained pragmatically by cross-validation.

Advanced training methods

Denoising auto-encoders

- Add to each pattern x a random noise and ask the auto-encoding network to reconstruct the original noise-free pattern x.
- This encourages the system to extract even stronger and more significant regularities from the input patterns

Advanced training methods(2)

Random dropout

- In stochastic backpropagation training, each hidden unit is randomly omitted from the network with probability 0.5.
- Each unit cannot rely on the presence of the other hidden units and is encouraged to identify useful information, independently of the other units.

Advanced training methods(3)

Curriculum learning

Training examples are presented to the network by starting from the easiest cases and then gradually proceeding to the more complex ones.

Gist

- Multi-layer perceptron neural networks (MLPs) are a flexible (non-parametric) modeling architecture composed of layers of sigmoidal units interconnected in a feedforward manner only between adjacent layers.
- Training from labeled examples can occur via variations of gradient descent (error backpropagation).
- Although gradient descent is a weak optimization method, it yields successful practical results.

Gist(2)

- Deep neural networks composed of many layers are becoming effective
- Learning schemes for deep MLP consist of:
- 1. an unsupervised preparatory phase
- 2. a final tuning phase using the scarce labeled examples.

Gist(3)

 To improve generalization, the use of controlled amounts of noise during training is effective

 Increasing the effort during training pays dividends in terms of improved generalization

Roberto Battiti • Mauro Brunato

The LION Way

Machine Learning plus Intelligent Optimization

The LOO Market and the April 2015

ROBERTO BATTITI, MAURO BRUNATO. *The LION Way: Machine Learning* plus *Intelligent Optimization*. LIONIab, University of Trento, Italy,

Apr 2015

http://intelligentoptimization.org/LIONbook



THE 11TH LEARNING AND INTELLIGENT OPTIMIZATION CONFERENCE http://intelligent-optimization.org/lion11/

DEC 10, 2016 PAPER SUBMISSION DEADLINE.

Thank you