

# Professional English – I

## Reading

Chixiao Chen

[cxchen@fudan.edu.cn](mailto:cxchen@fudan.edu.cn)

## Deep learning

Yann LeCun<sup>1,2</sup>, Yoshua Bengio<sup>3</sup> & Geoffrey Hinton<sup>4,5</sup>

Deep learning allows computational models that are composed of multiple processing layers to learn representations of data with multiple levels of abstraction. These methods have dramatically improved the state-of-the-art in speech recognition, visual object recognition, object detection and many other domains such as drug discovery and genomics. Deep learning discovers intricate structure in large data sets by using the backpropagation algorithm to indicate how a machine should change its internal parameters that are used to compute the representation in each layer from the representation in the previous layer. Deep convolutional nets have brought about breakthroughs in processing images, video, speech and audio, whereas recurrent nets have shone light on sequential data such as text and speech.

Machine-learning technology powers many aspects of modern society: from web searches to content filtering on social networks to recommendations on e-commerce websites, and it is increasingly present in consumer products such as cameras and smartphones. Machine-learning systems are used to identify objects in images, transcribe speech into text, match news items, posts or products with users' interests, and select relevant results of search. Increasingly, these applications make use of a class of techniques called deep learning.

intricate structures in high-dimensional data and is therefore applicable to many domains of science, business and government. In addition to beating records in image recognition<sup>1–4</sup> and speech recognition<sup>5–7</sup>, it has beaten other machine-learning techniques at predicting the activity of potential drug molecules<sup>8</sup>, analysing particle accelerator data<sup>9,10</sup>, reconstructing brain circuits<sup>11</sup>, and predicting the effects of mutations in non-coding DNA on gene expression and disease<sup>12,13</sup>. Perhaps more surprisingly, deep learning has produced extremely promising results for various tasks in natural language understanding<sup>14</sup>, particularly

# Authors

The ACM A.M. Turing Award is an annual prize given by the Association for Computing Machinery (ACM) to an individual selected for contributions "of lasting and major technical importance to the computer field". The Turing Award is generally recognized as the highest distinction in computer science and the "Nobel Prize of computing". The award is named after Alan Turing, a British mathematician and reader in mathematics at the University of Manchester. Turing is often credited as being the key founder of theoretical computer science and artificial intelligence. Since 2014, the award has been accompanied by a prize of US\$1 million, with financial support provided by Google.



Deep learning discovers **intricate** structure in large data sets by using the backpropagation algorithm to indicate how a machine should change its **internal** parameters that are used to compute the representation in each layer from the representation in the previous layer.

Deep **convolutional** nets have brought about **breakthroughs** in processing images, video, speech and audio, whereas recurrent nets have **shone** light on **sequential** data such as text and speech.

Machine-learning technology powers many aspects of modern society: from web searches to content filtering on social networks to recommendations on e-commerce websites, and it is increasingly present in consumer products such as cameras and smartphones.



For **decades**, constructing a pattern-**recognition** or machine-learning system required careful engineering and **considerable** domain expertise to design a feature extractor that transformed the **raw** data (such as the pixel values of an image) into a suitable internal representation or feature **vector** from which the learning **subsystem**, often a classifier, could detect or classify patterns in the input.

Deep-learning methods are representation-learning methods with multiple levels of representation, obtained by **composing** simple but **non-linear modules** that each transform the representation at one level (starting with the raw input) into a representation at a higher, slightly more **abstract** level.



An image, for example, comes in the form of an array of **pixel** values, and the learned features in the **first** layer of representation typically represent the **presence** or **absence** of **edges** at particular orientations and locations in the image. The **second** layer typically detects **motifs** by **spotting** particular arrangements of edges, **regardless** of small **variations** in the edge positions. The **third** layer may **assemble** motifs into larger combinations that correspond to parts of familiar objects, and **subsequent** layers would detect objects as combinations of these parts.

The key aspect of deep learning is that these layers of features are not designed by human engineers: they are learned from data using a general-purpose learning procedure.

In addition to beating **records** in image recognition [1]–[4] and speech recognition [5]–[7], it has beaten other machine-learning techniques at predicting the activity of potential drug **molecules** [8], analyzing **particle** accelerator data [9] [10] , **reconstructing** brain **circuits** [11], and predicting the effects of **mutations** in non-coding DNA on gene expression and disease [12],[13].

*Perhaps more surprisingly*, deep learning has produced extremely **promising** results for various tasks in natural language understanding, particularly topic classification, **sentiment** analysis, question answering and language translation.

These **adjustable** parameters, often called weights, are real numbers that can be seen as ‘**knobs**’ that define the input–output function of the machine.

The **objective** function, **averaged** over all the training examples, can be seen as a kind of hilly landscape in the high-dimensional space of weight values. The negative **gradient** vector indicates the direction of **steepest descent** in this landscape, taking it closer to a minimum, where the output error is low on average.

After training, the **performance** of the system is **measured** on a different set of examples called a test set. This serves to test the **generalization** ability of the machine — its ability to produce **sensible** answers on new inputs that it has never seen during training.



Problems such as image and speech recognition require the input–output **function** to be **insensitive** to **irrelevant** variations of the input, such as variations in **position**, **orientation** or **illumination** of an object, or variations in the **pitch** or **accent** of speech, while being very sensitive to particular **minute** variations (for example, the difference between a white wolf and a breed of wolf-like white dog called a Samoyed).

This is why **shallow** classifiers require a good feature extractor that solves the **selectivity–invariance dilemma** — one that produces representations that are selective to the aspects of the image that are important for discrimination, but that are invariant to irrelevant aspects such as the pose of the animal.

Trade-off

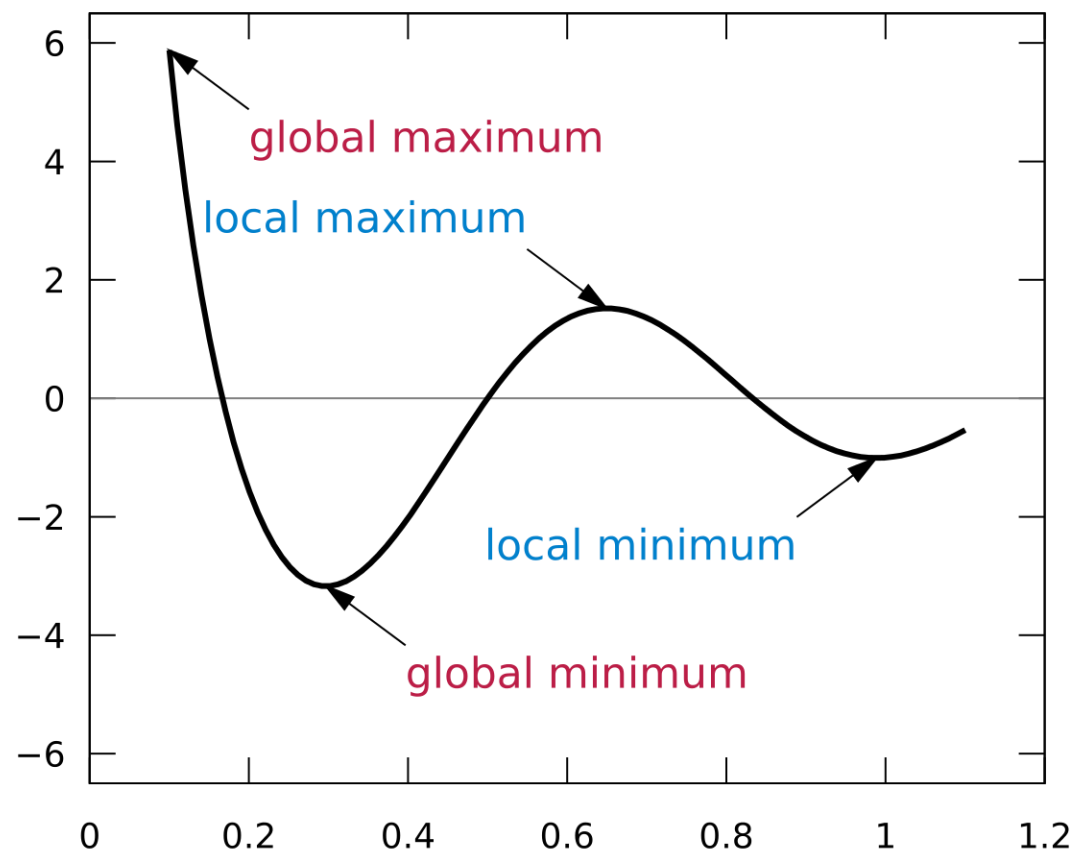
With multiple non-linear layers, say a **depth** of 5 to 20, a system can implement extremely intricate functions of its inputs that are **simultaneously** sensitive to minute details — **distinguishing** Samoyeds from white wolves — and insensitive to large irrelevant variations such as the background, **pose**, lighting and surrounding objects.

As long as the modules are relatively **smooth** functions of their inputs and of their internal weights, one can compute **gradients** using the backpropagation procedure.

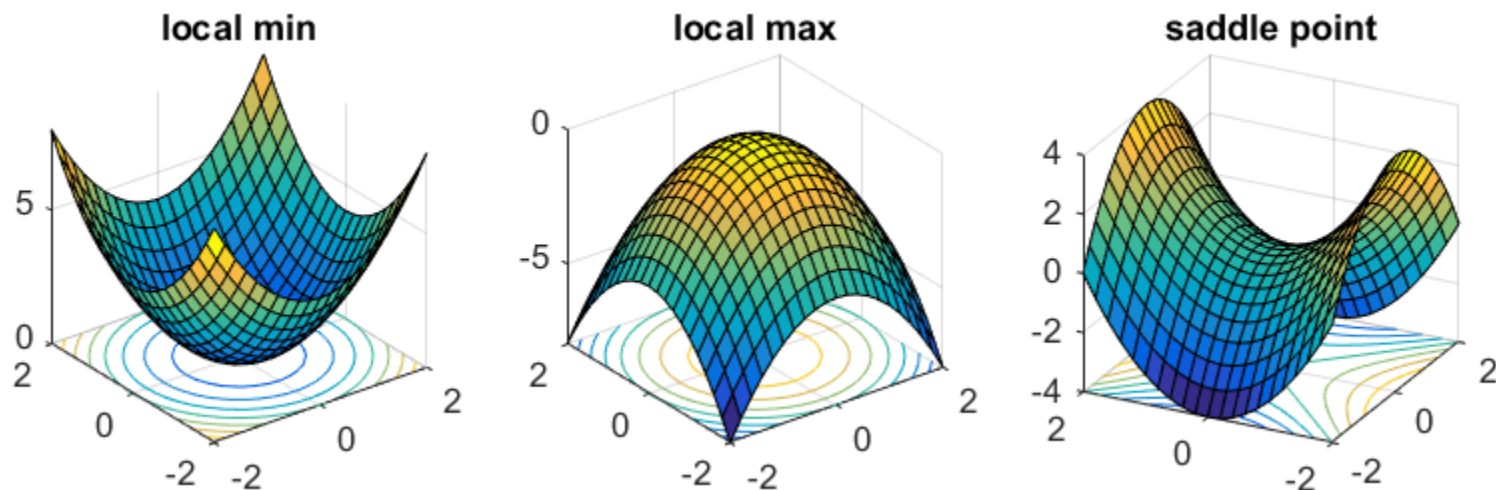
The backpropagation procedure to compute the gradient of an objective function with respect to the weights of a multilayer **stack** of modules is nothing more than a practical application of the chain rule for **derivatives**.

At present, the most popular non-linear function is the rectified linear unit (ReLU), which is simply the half-wave **rectifier**  $f(z) = \max(z, 0)$ . In past decades, neural nets used smoother non-linearities, such as  $\tanh(z)$  or  $1/(1 + \exp(-z))$ , but the ReLU typically learns much faster in networks with many layers, allowing training of a deep supervised network without unsupervised pre-training.

It was widely thought that learning useful, multistage, feature extractors with little **prior knowledge** was **infeasible**. In particular, it was commonly thought that simple gradient descent would get trapped in poor local **minima** — weight configurations for which no small change would reduce the average error.



Regardless of the **initial conditions**, the system nearly always reaches solutions of very similar **quality**. Recent **theoretical** and **empirical** results strongly suggest that local minima are not a serious issue in general. Instead, the landscape is packed with a combinatorially large number of **saddle points** where the gradient is zero, and the surface curves up in most dimensions and curves down in the remainders.





By ‘pre-training’ several layers of progressively more complex feature detectors using this reconstruction objective, the weights of a deep network could be initialized to sensible values.

Once deep learning had been rehabilitated, it turned out that the pre-training stage was only needed for small data sets.

In 2009, the approach was used to map short **temporal** windows of coefficients **extracted** from a sound wave to a set of **probabilities** for the various **fragments** of speech that might be represented by the frame in the center of the window. It achieved **record-breaking** results on a standard speech recognition **benchmark** that used a small **vocabulary** and was quickly developed to give record-breaking results on a large vocabulary task.

There was, however, one particular type of deep, feedforward network that was much easier to train and generalized much better than networks with full **connectivity** between **adjacent** layers. This was the convolutional neural network (ConvNet) .

Many data modalities are in the form of multiple arrays: 1D for signals and sequences, including language; 2D for images or audio spectrograms; and 3D for video or volumetric images.