

# ECBM E4040 Neural Networks and Deep Learning

## Lecture: Feedforward Deep Networks (cont'd)

Instructor Zoran Kostic

Columbia University  
Department of Electrical Engineering

Sep. 30th, 2016

# Part I

## Review of Previous Lecture

# Previous Lecture - Topics Covered

## Feedforward Deep Networks

- Multilayer Perceptrons from the 1980s
- Estimating Conditional Statistics
- Parametrizing a Learned Predictor
- Computational Graphs and Back-Propagation

# Previous Lecture - Learning Objectives

## Feedforward Deep Networks

- MLPs are realizations/implementations of compositions of multi-input and multi-output parametric functions.
- Examples of non-linear activation functions: sigmoid, tanh, softmax, radial basis functions, maxout, etc.
- Backpropagation: a method for computing gradients for multi-layer neural networks.
- The basic idea of the back-propagation algorithm is that the partial derivative of the cost  $J$  with respect to parameters  $\theta$  can be **decomposed recursively** by taking into consideration the composition of functions that relate  $\theta$  to  $J$ , via intermediate quantities that mediate that influence, e.g., the activations of hidden units in a deep neural network.

## Learning Objectives (cont'd)

### Back-Propagation

The back-propagation algorithm proceeds by first computing the gradient of the cost  $J$  with respect to output units, and these are used to compute the gradient of  $J$  with respect to the top hidden layer activations, which directly influence the outputs. The BPA continues computing the gradients of lower level hidden units one at a time in the same way.

The gradients on hidden and output units can be used to compute the gradient of  $J$  with respect to the parameters. In a typical network divided into many layers with each layer parameterized by a weight matrix and a vector of biases, the gradient on the weights and the biases is determined by the input to the layer and the gradient on the output of the layer.

## Part II

# Today's Lecture - Feedforward Deep Networks (cont'd)

# Implementing Stochastic Transformations

Traditional neural networks implement a deterministic transformation of some input variables  $\mathbf{x}$ . We can extend neural networks to implement stochastic transformations of  $\mathbf{x}$  by simply augmenting the neural network with randomly sampled inputs  $\mathbf{z}$ .

The neural network can then continue to perform deterministic computation internally, but the function  $f(\mathbf{x}, \mathbf{z})$  will appear stochastic to an observer who does not have access to  $\mathbf{z}$ . Provided that  $f$  is continuous and differentiable, we can then generally apply back-propagation as usual.

Let  $y$  be obtained by sampling the Gaussian distribution  $\mathcal{N}(\mu, \sigma)$ . We interpret the sampling process as the transformation

$$y = \mu + \sigma z,$$

where  $z$  is a Gaussian random variable  $\mathcal{N}(0, 1)$ .

## Implementing Stochastic Transformations (cont'd)

We are now able to back-propagate through the sampling operation, by regarding it as a deterministic operation with an additional input  $z$ .

Back-propagation through the sampling operation enables us to incorporate the transformation into a larger graph and, thereby, compute the derivatives of the loss function of interest  $J(y)$ .

We can also introduce functions that shape the distribution, e.g.,  $\mu = f(\mathbf{x}; \boldsymbol{\theta})$  and  $\sigma = g(\mathbf{x}; \boldsymbol{\theta})$  and use back-propagation through this functions to derive  $\nabla_{\boldsymbol{\theta}} J(y)$ .



# Implementing Stochastic Transformations (cont'd)

More generally, using the representation

$$p(y|x; \theta) = p(y|\omega),$$

where  $\omega$  is a variable containing both parameters  $\theta$  and the inputs  $x$ , we have

$$y = f(z, \omega),$$

where  $z$  is a source of randomness. Here,  $\omega$  must not be a function of  $z$ , and  $z$  must not be a function of  $\omega$ . This is often called the **reparametrization trick**.

In neural network applications,  $z$  is typically drawn from some simple distribution, such as a unit uniform or unit Gaussian distribution. Complex distributions are achieved by allowing the deterministic portion of the network to reshape its input.

# The Universal Approximation Theorem

## The Power and Limitations of the UAP

### Theorem (Universal Approximation Theorem)

*A feedforward network with a linear output layer and at least one hidden layer with any squashing activation function (such as the logistic sigmoid activation function) can approximate any Borel measurable function from one finite-dimensional space to another with any desired non-zero amount of error, provided that the network is given enough hidden units.*

The derivatives of the feedforward network can also approximate the derivatives of the function arbitrarily well.

# The Universal Approximation Theorem (cont'd)

## The Power and Limitations of the UAP

Any continuous function on a closed and bounded subset of  $\mathbb{R}^n$  is Borel measurable and therefore may be approximated by a neural network.

A neural network may also approximate any function mapping from any finite dimensional discrete space to another. Interestingly, universal approximation theorems have also been proven for a wider class of non-linearities which includes the now commonly used rectified linear unit.

# The Universal Approximation Theorem (cont'd)

## The Power and Limitations of the UAP

Feedforward networks provide a universal system for representing functions: given a function, we can find a large MLP that approximates/represents the function.

Even if the MLP is able to represent a function, there are no guarantees that the training algorithm will be able to learn it.

Learning can fail for two different reasons:

- the optimization algorithm used for training may not be able to find the value of the parameters that corresponds to the desired function;
- the training algorithm might choose the wrong function due to overfitting.

There is no universal procedure for examining a training set of specific examples and choosing a function that will generalize to points not in the training set.

# The Universal Approximation Theorem (cont'd)

## The Power and Limitations of the UAP

The UAT says that there exists a network large enough to achieve any degree of accuracy; it does not say, however, how large this network will be.

Unfortunately, in the worse case, an exponential number of hidden units may be required. This is easiest to see in the binary case: the number of possible binary functions on vectors  $v \in \{0, 1\}^n$  is  $2^{2^n}$  and selecting one such function requires  $2^n$  bits, which will in general require  $\mathcal{O}(2^n)$  degrees of freedom.

### Remark

*A feedforward network with a single layer is sufficient to represent any function, but it may be infeasibly large and may fail to learn and generalize correctly. Both of these failure modes suggest the use of deeper models.*

# The Universal Approximation Theorem (cont'd)

## The Power and Limitations of the UAP

From a representation learning point of view the learning with a deep architecture

- consists of discovering a set of underlying factors of variation that can in turn be described in terms of other, simpler underlying factors of variation;
- expresses the belief that the function we want to learn is a computer program consisting of multiple steps, where each step makes use of the previous steps output. These intermediate outputs are not necessarily factors of variation, but can instead be analogous to counters or pointers that the network uses to organize its internal processing.

Using deep architectures does indeed express a useful prior over the space of functions the model learns.

# Improved Performance with Piecewise Linear Hidden Units

Recent performance improvements of deep neural networks can be attributed to

- **increases in computational power;**
- **the size of datasets.**

The main algorithmic improvements are due to

- the **use of piecewise linear units**, such as absolute value rectifiers and rectified linear units. Such units consist of two linear pieces and their behavior is driven by a single weight vector.
- using a rectifying non-linearity is the single most important factor in improving the performance of a recognition system among several different factors of neural network architecture design.

## Improved Performance with Piecewise LHUs (cont'd)

For small datasets, **using rectifying nonlinearities is even more important than learning the weights** of the hidden layers. Random weights are sufficient to propagate useful information through a rectified linear network, allowing the classifier layer at the top to learn how to map different feature vectors to class identities.

Learning is far easier in deep rectified linear networks than in deep networks that have curvature or two-sided saturation in their activation functions.

Just as piecewise linear networks are good at propagating information forward, back-propagation in such a network is also piecewise linear and propagates information about the error derivatives to all of the gradients in the network.



## Improved Performance with Piecewise LHUs (cont'd)

The use rectified linear units is strongly motivated by biological considerations. The half-rectifying non-linearity was intended to capture the following properties of biological neurons:

- for some inputs, biological neurons are completely inactive.
- for some inputs, a biological neurons output is proportional to its input.
- most of the time, biological neurons operate in the regime where they are inactive (e.g., they have sparse activations).

### Remark

*Piecewise linear units represent by far the most popular class of hidden units.*