

# Multi-layer auto-encoders

a CS221 artificial intelligence project

Charles Celerier	Bill Chickering
Computer Science	Computer Science
Stanford	Stanford
cceleri@cs.stanford.edu	chickering@cs.stanford.edu

August 18, 2013

This is a short introduction.

## 1 Introduction

A fundamental challenge to supervised machine learning is that of feature selection. There exist many strategies for determining a set of variables derived from an original representation to accurately classify domain data. Ideally, the selected variables minimize redundant and/or irrelevant information and thereby reduce the risk of overfitting. Historically, domain experts and computer scientists have laboriously handcrafted feature sets for a variety of machine learning problems. These manually engineered features suffer from the requirements of domain expertise and a great deal of work, and therefore, might not scale well. Moreover, such feature sets can inadvertently ignore important structure in the data. These deficiencies can be addressed via unsupervised feature learning using dimension reduction techniques such as principal component analysis (PCA) or factor analysis. Recently it was found that K-means clustering offers another straightforward method for determining an efficacious feature space (e.g. the CS221 Visual Cortex programming assignment). In this report we document a project in which unsupervised feature learning is performed using autoencoder neural networks to discover feature spaces that can be utilized by a supervised learning algorithm to effectively classify images of handwritten digits.

Our subject is the MNIST database of handwritten digits<sup>1</sup>, which includes 50,000 training and 10,000 testing images, each a 28x28 grid of real numbers between 0 and 1. It is worth mentioning that this dataset has been employed before to demonstrate feature learning with autoencoders. Indeed, the pioneering paper by Hinton *et al.* [2] that introduced the technique upon which our project is based used this dataset among others. We chose to work with the MNIST data because of its large size, simple format, and obvious classifications thereby allowing us to focus on technique.

## 2 Method

An autoencoder is a special type of neural network consisting of an input layer, an output layer, and one or more hidden layers. An autoencoder aims to approximate the identity function as closely as possible, reproducing the signals of its input layer at its output layer. By placing constraints on the network, input signals must be compressed and decompressed using optimized encodings in order to approximately reproduce the input. Such encodings can then serve as feature sets in a supervised learning problem. The challenge is training such autoencoders in order to discover the optimal encoding for a particular domain space.

Until recently, the standard approach to training multilayer feedforward neural networks, such as autoencoders, is to use a variation of the backpropagation algorithm. We describe this algorithm later in this report. For now, we merely point out a key feature of the backpropagation algorithm: it attempts to optimize all weights and biases of a neural network simultaneously. In the case of large networks containing multiple layers and 1000+ nodes, backpropagation following random initialization of the parameters can fail due to

---

<sup>1</sup><http://yann.lecun.com/exdb/mnist/>

local minima in the objective function [3]. A technique pioneered by Hinton *et al.* [2] largely overcomes this challenge by “pretraining” individual restricted Boltzmann machines (RBMs), which can then be “unrolled” to form a deep, multilayer autoencoder with weights and biases that are close to optimal prior to running a backpropagation algorithm.

An RBM may be represented as a complete bipartite graph where one set corresponds to “visible” units  $\nu_i$  and the other to “hidden” units  $h_j$ . We work with binary RBMs in which each node is either on or off. A weight  $w_{ij}$  is associated with each edge that connects a visible unit  $\nu_i$  to a hidden unit  $h_j$ . Given the state of all visible units, the probability that a hidden unit is in the “on” state is given by  $\sigma(b_j + \sum_i \nu_i w_{ij})$ , where  $\sigma(x)$  is the logistic function  $1/(1 + \exp(-x))$  and  $b_j$  is the bias associated with hidden unit  $j$ . By coupling the visible and hidden units in this way, an RBM models a particular statistical ensemble of unit configurations, which in our case are images. The technique of Hinton *et al.* “pretrains” an RBM so that it approximately models a training set.

A detailed description of the technique of Hinton *et al.* is beyond the scope of this report; we refer the interested reader to the original papers referenced here as well as to our provided source code. The basic idea is to model our data as a physical system. We may then borrow some statistical machinery from physics to help pretrain our RBMs. Perhaps the most fundamental of these is the *free energy* ascribed to each training image, given by

$$E(\boldsymbol{\nu}, \mathbf{h}) = - \sum_{i \in \text{visible}} b_i \nu_i - \sum_{j \in \text{hidden}} b_j h_j - \sum_{i,j} \nu_i h_j w_{ij}. \quad (1)$$

As per the Maxwell-Boltzmann distribution, the probability of a particular image is proportional to  $\exp(-E(\boldsymbol{\nu}, \mathbf{h}))$ . Using stochastic gradient descent to maximize the log-likelihood of our training images, along with an approximation technique known as *contrastive divergence* that allows us estimate an expectation over all possible unit configurations under the probability distribution modeled by an RBM, we iteratively update the biases and weights of an RBM. Once an RBM is optimized in this manner, it is straightforward to create an autoencoder with a single hidden layer by adding an output layer that simply decodes the states of the hidden layer. To construct an autoencoder with three hidden layers, we pretrain an additional RBM using the output (i.e. hidden layer) of the first RBM as its visible layer. Importantly, the hidden layer of the second RBM must consist of fewer nodes than the hidden layer of the first RBM. In this way, multiple RBMs, each with fewer hidden nodes than the one before it, can be unrolled to form a pretrained deep neural network autoencoder.

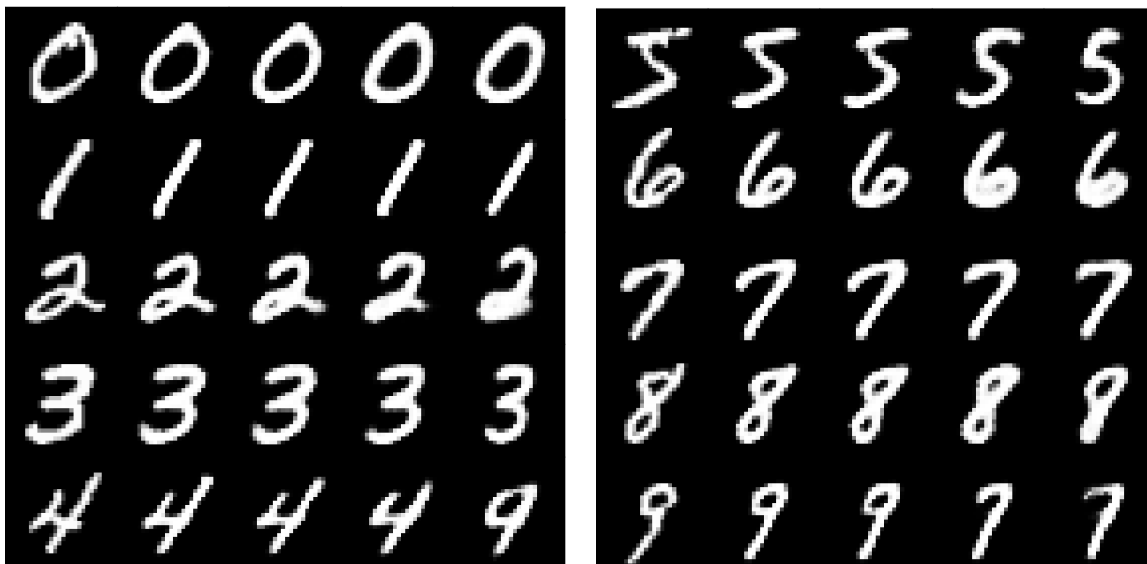
Once pretrained, an autoencoder can be fine-tuned using backpropagation without fear of local minima in our objective function (i.e. the mean squared error). As with the RBM pretraining technique, the details of this classic algorithm for training neural networks are beyond this report’s scope and are well documented elsewhere (an example implementation is also provided within our source code). The essence of this algorithm is to propagate the discrepancy between the autoencoder’s output and its input back through the neural network, adjusting its parameters along the way.

These autoencoders essentially function by compressing and decompressing the input image. The compression algorithm used has been optimized for the given training data. In addition to providing a domain-optimized compression algorithm, we can also use the dimensionally reduced representation of the data as a feature space. By removing the decoders such that the innermost hidden layer with the fewest nodes becomes the output layer, we construct a feature encoder. To explore the utility of such learned feature spaces, we conducted several supervised learning experiments in which we construct naïve Bayes classifiers and compare their error rates.

### 3 Results

During the course of our project, we constructed several autoencoders consisting of a variety of dimensions. We focus here on autoencoders constructed from RBMs consisting of 3000, 1500, 750, 100, 30, and 10 node hidden layers. In addition to these hidden layers, all autoencoders have 784 node input and output layers. For example, our narrowest waist autoencoder consists of 784-3000-1500-750-100-30-10-30-100-750-1500-3000-784 node layers, while our fattest waist autoencoder consists of 784-3000-784 node layers. In Fig. 1 we show several autoencoded handwritten digits. From the leftmost to rightmost columns we have: the original image, encoded using 784-3000-1500-750, 784-3000-1500-750-100, 784-3000-1500-750-100-30, and 784-3000-1500-750-100-30-10. As expected, the quality of the reproduction reduces from left to right.

To test the utility of the learned feature spaces within the context of supervised learning, we trained several naive Bayes classifiers. The error rates obtained for both the training and test data are shown in Fig.



(a) Encoding of digits 0 – 4 in each layer

(b) Encoding of digits 5 – 9 in each layer

Figure 1: Intermediate encodings for 784-...-750-100-30-10 multilayer autoencoder.

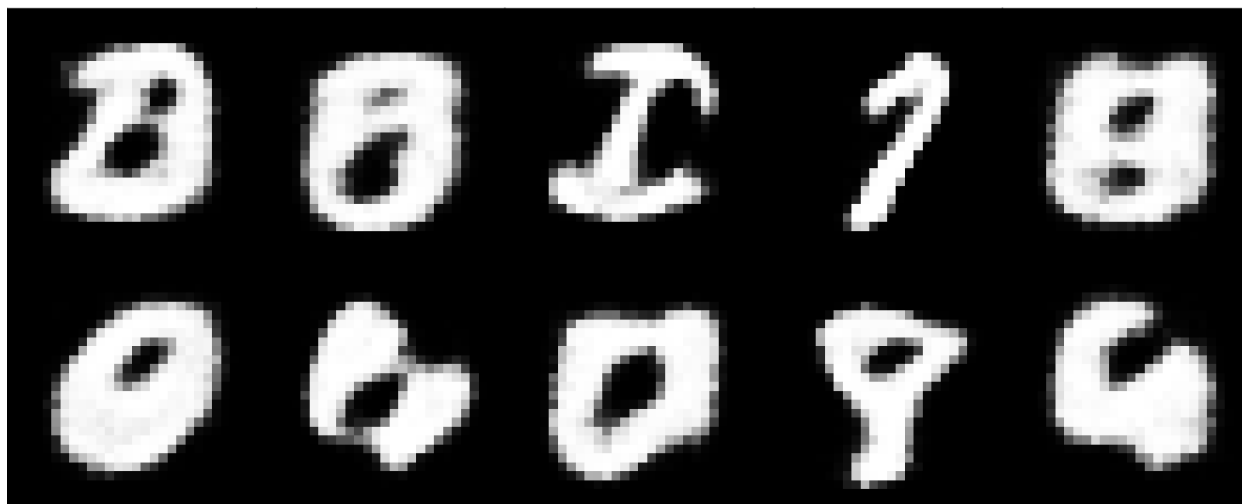


Figure 2: Highest level detectors for 784-3000-1500-750-100-30-10.

3. Column A reveals a training error of about 17% and test error of about 15% using the raw data. This is our baseline. In column B we see the error rates are essentially unchanged using the expanded 3000 feature nodes. Remarkably, for columns C and D, corresponding to 1500 and 750 feature nodes, respectively, we small but significant reductions in the error rate. More striking, however, is that column E, which corresponds to 100 feature nodes, also shows a reduced error rate despite a dimensionally reduced representation of the data. As we continue to reduce the dimensionality of our feature space, columns F and G, the error rates increase.

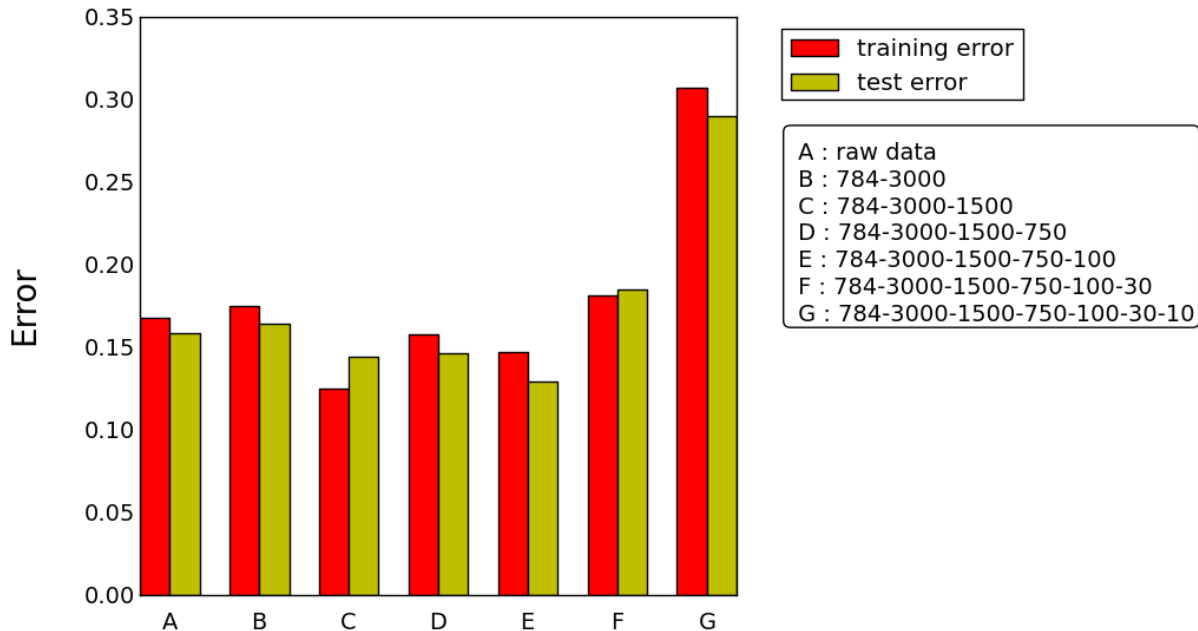


Figure 3: Training and test error of NB classifier on highest-level autoencoder features

## 4 Conclusion

## References

- [1] *Restricted Boltzmann Machines (RBM)*, <http://deeplearning.net/tutorial/rbm.html>.
- [2] G. E. Hinton and R. R. Salakhutdinov, *Reducing the dimensionality of data with neural networks*, *Science* **313** (2006), no. 5786, 504–507.
- [3] N. Kambhatla and T. Leen., *Dimension reduction by local principal component analysis*, *Neural Computation* **9** (1997), no. 7, 1493–1516.
- [4] Raul Rojas, *Neural networks - a systematic introduction*, <http://page.mi.fu-berlin.de/rojas/neural/chapter/K7.pdf>, 1996.