

Autoencoders, Unsupervised Learning, and Deep Architectures

Pierre Baldi

PFBALDI@ICS.UCI.EDU

*Department of Computer Science
University of California, Irvine
Irvine, CA 92697-3435*

Editor: I. Guyon, G. Dror, V. Lemaire, G. Taylor and D. Silver

Abstract

Autoencoders play a fundamental role in unsupervised learning and in deep architectures for transfer learning and other tasks. In spite of their fundamental role, only linear autoencoders over the real numbers have been solved analytically. Here we present a general mathematical framework for the study of both linear and non-linear autoencoders. The framework allows one to derive an analytical treatment for the most non-linear autoencoder, the Boolean autoencoder. Learning in the Boolean autoencoder is equivalent to a clustering problem that can be solved in polynomial time when the number of clusters is small and becomes NP complete when the number of clusters is large. The framework sheds light on the different kinds of autoencoders, their learning complexity, their horizontal and vertical composability in deep architectures, their critical points, and their fundamental connections to clustering, Hebbian learning, and information theory.

Keywords: autoencoders, unsupervised learning, compression, clustering, principal component analysis, boolean, complexity, deep architectures, hebbian learning, information theory

1. Introduction

Autoencoders are simple learning circuits which aim to transform inputs into outputs with the least possible amount of distortion. While conceptually simple, they play an important role in machine learning. Autoencoders were first introduced in the 1980s by Hinton and the PDP group (Rumelhart et al., 1986) to address the problem of “backpropagation without a teacher”, by using the input data as the teacher. Together with Hebbian learning rules (Hebb, 1949; Oja, 1982), autoencoders provide one of the fundamental paradigms for unsupervised learning and for beginning to address the mystery of how synaptic changes induced by local biochemical events can be coordinated in a self-organized manner to produce global learning and intelligent behavior.

More recently, autoencoders have taken center stage again in the “deep architecture” approach (Hinton et al., 2006; Hinton and Salakhutdinov, 2006; Bengio and LeCun, 2007; Erhan et al., 2010) where autoencoders, particularly in the form of Restricted Boltzmann Machines (RBMS), are stacked and trained bottom up in unsupervised fashion, followed by a supervised learning phase to train the top layer and fine-tune the entire architecture. The bottom up phase is agnostic with respect to the final task and thus can obviously be

used in transfer learning approaches. These deep architectures have been shown to lead to state-of-the-art results on a number of challenging classification and regression problems.

In spite of the interest they have generated, and with a few exceptions (Baldi and Hornik, 1988; Sutskever and Hinton, 2008; Montufar and Ay, 2011), little theoretical understanding of autoencoders and deep architectures has been obtained to this date. Additional confusion may have been created by the use of the term “deep”. A deep architecture from a computer science perspective should have n^α polynomial-size layers, for some small $\alpha > 0$, where n is the size of the input vectors (see Clote and Kranakis (2002) and references therein). But that is not the case in the architectures described in Hinton et al. (2006) and Hinton and Salakhutdinov (2006), which seem to have constant or at best logarithmic depth, the distinction between finite and logarithmic depth being almost impossible for the typical values of n used in computer vision, speech recognition, and other typical problems. Thus the main motivation behind this work is to derive a better theoretical understanding of autoencoders, with the hope of gaining better insights into the nature of unsupervised learning and deep architectures.

If general theoretical results about deep architectures exist, these are unlikely to depend on a particular hardware realization, such as RBMs. Similar results ought to be true for alternative, or more general, forms of computation. Thus the strategy proposed here is to introduce a general framework and study different kinds of autoencoder circuits, in particular Boolean autoencoders which can be viewed as the most extreme form of non-linear autoencoders. The expectation is that certain properties of autoencoders and deep architectures may be easier to identify and understand mathematically in simpler hardware embodiments, and that the study of different kinds of autoencoders may facilitate abstraction and generalization through the identifications of common properties.

For this purpose, we begin in Section 2 by describing a fairly general framework for studying autoencoders. In Section 3, we review and extend the known results on linear autoencoders. In the light of deep architectures, we look at novel properties such a vertical composition (stacking) and connection of critical points to stability under recycling (feeding outputs back to the input layer). In Section 4, we study Boolean autoencoders, and prove several properties including their fundamental connection to clustering. In Section 5, we address the complexity of Boolean autoencoder learning. In Section 6, we study autoencoders with large hidden layers, and introduce the notion of horizontal composition of autoencoders. In Section 7, we address other classes of autoencoders and generalizations. Finally, in Section 8, we summarize the results and their possible consequences for the theory of deep architectures.

2. A General Autoencoder Framework

To derive a fairly general framework, an $n/p/n$ autoencoder (Figure 1) is defined by a t-tuple $n, p, m, \mathbb{F}, \mathbb{G}, \mathcal{A}, \mathcal{B}, \mathcal{X}, \Delta$ where:

1. \mathbb{F} and \mathbb{G} are sets.
2. n and p are positive integers. Here we consider primarily the case where $0 < p < n$.
3. \mathcal{A} is a class of functions from \mathbb{G}^p to \mathbb{F}^n .

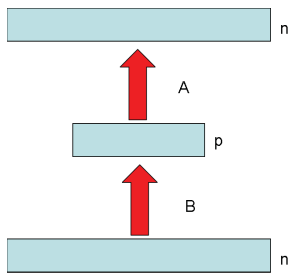


Figure 1: An $n/p/n$ Autoencoder Architecture.

4. \mathcal{B} is a class of functions from \mathbb{F}^n to \mathbb{G}^p .
5. $\mathcal{X} = \{x_1, \dots, x_m\}$ is a set of m (training) vectors in \mathbb{F}^n . When external targets are present, we let $\mathcal{Y} = \{y_1, \dots, y_m\}$ denote the corresponding set of target vectors in \mathbb{F}^n .
6. Δ is a dissimilarity or distortion function (e.g. L_p norm, Hamming distance) defined over \mathbb{F}^n .

For any $A \in \mathcal{A}$ and $B \in \mathcal{B}$, the autoencoder transforms an input vector $x \in \mathbb{F}^n$ into an output vector $A \circ B(x) \in \mathbb{F}^n$ (Figure 1). The corresponding *autoencoder problem* is to find $A \in \mathcal{A}$ and $B \in \mathcal{B}$ that minimize the overall distortion function:

$$\min E(A, B) = \min_{A, B} \sum_{t=1}^m E(x_t) = \min_{A, B} \sum_{t=1}^m \Delta(A \circ B(x_t), x_t) \quad (1)$$

In the non auto-associative case, when external targets y_t are provided, the minimization problem becomes:

$$\min E(A, B) = \min_{A, B} \sum_{t=1}^m E(x_t, y_t) = \min_{A, B} \sum_{t=1}^m \Delta(A \circ B(x_t), y_t) \quad (2)$$

Note that $p < n$ corresponds to the regime where the autoencoder tries to implement some form of compression or feature extraction. The case $p \geq n$ is discussed towards the end of the paper.

Obviously, from this general framework, different kinds of autoencoders can be derived depending, for instance, on the choice of sets \mathbb{F} and \mathbb{G} , transformation classes \mathcal{A} and \mathcal{B} , distortion function Δ , as well as the presence of additional constraints, such as regularization. To the best of our knowledge, neural network autoencoders were first introduced by the PDP group as a special case of this definition, with all vectors components in $\mathbb{F} = \mathbb{G} = \mathbb{R}$ and A and B corresponding to matrix multiplications followed by non-linear sigmoidal transformations with an L_2^2 error function. [For regression problems, the non-linear sigmoidal transformation is typically used only in the hidden layer]. As an approximation to this case, in the next section, we study the linear case with $\mathbb{F} = \mathbb{G} = \mathbb{R}$. More generally, linear autoencoders correspond to the case where \mathbb{F} and \mathbb{G} are fields and \mathcal{A} and \mathcal{B} are the classes of linear transformations, hence A and B are matrices of size $p \times n$ and $n \times p$ respectively. The

linear real case where $\mathbb{F} = \mathbb{G} = \mathbb{R}$ and Δ is the squared Euclidean distance was addressed in Baldi and Hornik (1988) (see also Broulard and Kamp (1988)). More recently the theory has been extended also to complex-valued linear autoencoders (Baldi et al., 2011)

3. The Linear Autoencoder

We partly restate without proof the results derived in Baldi and Hornik (1988) for the linear real case with $\Delta = L_2^2$, but organize them in a way meant to highlight the connections to other kinds of autoencoders, and extend their results from a deep architecture perspective. We use A^t to denote the transpose of any matrix A .

1) Group Invariance. Every solution is defined up to multiplication by an invertible $p \times p$ matrix C , or equivalently up to a change of coordinates in the hidden layer. This is obvious since since $AC^{-1}CB = AB$.

2) Problem Complexity. While the cost function is quadratic and all the operations are linear, the overall problem is not convex because the hidden layer limits the rank of the overall transformation to be at most p , and the set of matrices of rank p or less is *not* convex. However the linear autoencoder problem over \mathbb{R} can be solved analytically. Note that in this case one is interested in finding a low rank approximation to the identity function.

3) Fixed Layer Solution. The problem becomes convex if A is fixed, or if B is fixed. When A is fixed, assuming A has rank p and that the data covariance matrix $\Sigma_{XX} = \sum_i x_i x_i^t$ is invertible, then at the optimum $B^* = B(A) = (A^t A)^{-1} A^t$. When B is fixed, assuming B has rank p and that Σ_{XX} is invertible, then at the optimum $A^* = A(B) = \Sigma_{XX} B^t (B \Sigma_{XX} B^t)^{-1}$.

4) The Landscape of E . The overall landscape of E has no local minima. All the critical points where the gradient of E is zero, correspond to projections onto subspaces associated with subsets of eigenvectors of the covariance matrix Σ_{XX} . Projections onto the subspace associated with the p largest eigenvalues correspond to the global minimum and Principal Component Analysis. All other critical point, corresponding to projections onto subspaces associated with other set of eigenvalues, are saddle points. More precisely, if $\mathcal{I} = i_1, \dots, i_p$ ($1 \leq i_1 < \dots < i_p \leq n$) is any ordered list of indices, let $U_{\mathcal{I}} = [u_1, \dots, u_p]$ denote the matrix formed by the orthonormal eigenvectors of Σ_{XX} associated with the eigenvalues $\lambda_{i_1}, \dots, \lambda_{i_p}$. Then two matrices A and B of rank p define a critical point if and only if there is a set \mathcal{I} and an invertible $p \times p$ matrix C such that $A = U_{\mathcal{I}} C$, $B = C^{-1} U_{\mathcal{I}}^t$, and $W = AB = P_{U_{\mathcal{I}}}$, where $P_{U_{\mathcal{I}}}$ is the orthogonal projection onto the subspace spanned by the columns of $U_{\mathcal{I}}$. At the global minimum, assuming that $C = I$, the activities in the hidden layer are given by the dot products $u_1^t x \dots u_p^t x$ and correspond to the coordinates of x along the first p eigenvectors of Σ_{XX} .

5) Clustering. The global minimum performs a form of clustering by hyperplane, with respect to $\text{Ker} B$, the kernel of B . For any given vector x , all the vectors of the form $x + \text{Ker}(B)$ are mapped onto the same vector $y = AB(x) = AB(x + \text{Ker} B)$.

6) Recycling Stability. At any critical point, AB is a projection operator and thus recycling outputs is stable at the first pass: $(AB)^n(x) = AB(x) = U_{\mathcal{I}} U_{\mathcal{I}}^t(x)$ for any $n \geq 1$.

7) Generalization. At any critical point, for any x , $AB(x)$ is equal to the projection of x onto the corresponding subspace and the corresponding error can be expressed easily as the squared distance of x to the projection space.

8) Vertical Composition. The global minimum of E remains the same if additional

matrices of rank greater or equal to p are introduced between the input layer and the hidden layer or the hidden layer and the output layer. Thus there is no reduction in overall distortion by introducing such matrices. However, if such matrices are introduced for other reasons, there is a composition law so that the optimum solution for a deep autoencoder with a stack of matrices, can be obtained by combining the optimal solutions of shallow autoencoders. More precisely, consider an autoencoder network with layers of size $n/p_1/p/p_1/n$ (Figure 2) with $n > p_1 > p$. Then the optimal solution of this network can be obtained by first computing the optimal solution for an $n/p_1/n$ autoencoder network, and combining it with the optimal solution of an $p_1/p/p_1$ autoencoder network using the activities in the hidden layer of the first network as the training set for the second network, exactly as in the case of stacked RBMs (Hinton et al., 2006; Hinton and Salakhutdinov, 2006). This is because the projection onto the subspace spanned by the top p eigenvectors can be composed by a projection onto the subspace spanned by the top p_1 eigenvectors, followed by a projection onto the subspace spanned by the top p eigenvectors.

9) External Targets. With the proper adjustments, the results above remain essentially the same if a set of target output vectors y_1, \dots, y_m is provided, instead of x_1, \dots, x_m serving as the targets (see (Baldi and Hornik, 1988)).

10) Symmetries and Hebbian Rules. At the global minimum, for $C = I$, $A = B^t$. The constraint $A = B^t$ can be imposed during learning by “weight sharing” and is consistent with a Hebbian rule that is symmetric between the pre- and post synaptic neurons and is applied to the network by clamping the output units to be equal to the input units (or having a folded autoencoder).

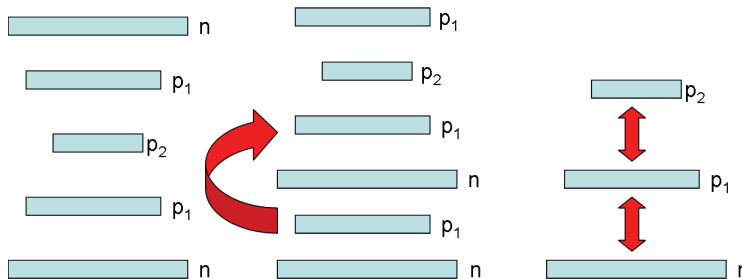


Figure 2: Vertical Composition of Autoencoders.

4. The Boolean Autoencoder

The Boolean autoencoder is the most extreme form of non-linear autoencoder. In the purely Boolean case, we have $\mathbb{F} = \mathbb{G} = \{0, 1\}$, A and B are unrestricted Boolean functions, and Δ is the Hamming distance. Many variants of this problem can be obtained by restricting the classes \mathcal{A} and \mathcal{B} of Boolean functions, for instance by bounding the connectivity of the hidden units. The linear case with $\mathbb{F} = \mathbb{G} = \{0, 1\} = \mathbb{F}_2$, where \mathbb{F}_2 is the Galois field with two elements, is a special case of the Boolean case and will be discussed later. For lack of space, proofs are only briefly sketched.

1) Group Invariance. Every solution is defined up to a permutation of the 2^p points of the hypercube \mathbb{H}^p . This is because the Boolean functions are unrestricted and therefore their lookup tables can accommodate any such permutation, or relabeling of the hidden states.

2) Problem Complexity. In general, the overall optimization problem is NP-hard. To be more precise, one must specify the regime of interest characterized by which variables among n , m , and p are going to infinity. Obviously one must have $n \rightarrow \infty$. If p does not go to infinity, then the problem can be polynomial, for instance when the centroids must belong to the training set. If $p \rightarrow \infty$ and m is a polynomial in n , which is the case of interest in machine learning where typically m is a low degree polynomial in n , then the problem of finding the best Boolean mapping (i.e. the Boolean mapping that minimizes the distortion E associated with the Hamming distance on the training set) is NP hard, or the corresponding decision problem is NP-complete. More precisely the optimisation problem is NP-hard in the regime where $p \sim \epsilon \log_2 m$ with $\epsilon > 0$. A proof of this result is given in the next section.

3) Fixed Layer Solution. If the A mapping is fixed, then it is easy to find the optimal B mapping. Conversely if the B mapping is fixed, it is easy to find the optimal A mapping. To see this, assume first that A is fixed. Then for each of the 2^p possible Boolean vectors h_1, \dots, h_{2^p} of the hidden layer, $A(h_1), \dots, A(h_{2^p})$ provide 2^p points (centroids) in the hypercube \mathbb{H}^n . One can build the corresponding Voronoi partition by assigning each point of \mathbb{H}^n to its closest centroid, breaking ties arbitrarily, thus forming a partition of \mathbb{H}^n into 2^p corresponding clusters $\mathcal{C}_1, \dots, \mathcal{C}_{2^p}$, with $\mathcal{C}_i = \mathcal{C}^{Vor}(A(h_i))$. The optimal mapping B^* is then easily defined by setting $B^*(x) = h_i$ for any x in $\mathcal{C}_i = \mathcal{C}^{Vor}(A(h_i))$. Conversely, assume that B is fixed. Then for each of the 2^p possible Boolean vectors h_1, \dots, h_{2^p} of the hidden layer, let $\mathcal{C}^B(h_i) = \{x \in \mathbb{H}^n : B(x) = h_i\}$. To minimize the reconstruction error, A^* must map h_i onto a point y of \mathbb{H}^n minimizing the sum of Hamming distances to points in $\mathcal{X} \cap \mathcal{C}^B(h_i)$. It is easy to see that the minimum is realized by the component-wise majority vector $A^*(h_i) = \text{Majority}[\mathcal{X} \cap \mathcal{C}^B(h_i)]$, breaking ties arbitrarily (e.g. by a coin flip).

4) The Landscape of E. In general E has many local minima (e.g. with respect to the Hamming distance applied to the lookup tables of A and B). Critical points are defined to be the points satisfying simultaneously the equations above for A^* and B^* .

5) Clustering. The overall optimization problem is a problem of optimal clustering. The clustering is defined by the transformation B . Approximate solutions can be sought by many algorithms, such as k-means, belief propagation (Frey and Dueck, 2007), minimum spanning paths and trees (Slagle et al., 1975), and hierarchical clustering.

6) Recycling Stability. At any critical point, recycling outputs is stable at the first pass so that for any x $(AB)^n(x) = AB(x)$ (and is equal to the majority vector of the corresponding Voronoi cluster).

7) Generalization. At any critical point, for any x , $AB(x)$ is equal to the centroid of the corresponding Voronoi cluster and the corresponding error can be expressed easily.

8) Vertical Composition. The global optimum remains the same if additional Boolean layers of size equal or greater to p are introduced between the input layer and the hidden layer and/or the hidden layer and the output layer. Thus there is no reduction in overall distortion E by adding such layers. Consider a Boolean autoencoder network with layers of size n, p_1, p, p_1, n (Figure 2) with $n > p_1 > p$. Then the optimal solution of this network can be obtained by first computing the optimal solution for an n, p_1, n autoencoder network,

and combining it with the optimal solution of an p_1, p, p_1 autoencoder network using the activity in the hidden layer of the first network as the training set, exactly as in the case of stacked RBMs. The reason for this is that the global optimum correspond to clustering into 2^p clusters, and this can be obtained by first clustering into 2^{p_1} clusters, and then clustering these clusters into 2^p clusters. The stack of Boolean functions performs hierarchical clustering with respect to the input space.

9) External Targets. With the proper adjustments, the results above remain essentially the same if a set of target output vectors y_1, \dots, y_m is provided, instead of x_1, \dots, x_m serving as the targets. To see this, consider a deep architecture consisting of a stack of autoencoders along the lines of [Hinton et al. \(2006\)](#). For any activity vector h in the last hidden layer before the output layer, compute the set of points $\mathcal{C}(h)$ in the training set that are mapped to h by the stacked architecture. Assume, without any loss of generality, that $\mathcal{C}(h) = \{x_1, \dots, x_k\}$ with corresponding targets $\{y_1, \dots, y_k\}$. Then it is easy to see that the final output for h produced by the top layer ought to be the centroid of the targets given by $Majority(y_1, \dots, y_k)$

5. Clustering Complexity on the Hypercube

In this section, we briefly review some results on clustering complexity and then prove that the hypercube clustering decision problem is in general NP-complete. The complexity of various clustering problems, in different spaces, or with different objective functions, has been studied in the literature. There are primarily two kind of results: (1) graphical results derived on graphs $G = (V, E, \Delta)$ where the dissimilarity Δ is not necessarily a distance; and (2) geometric results derived in the Euclidean space \mathbb{R}^d where $\Delta = L_2^2, L_2$, or L_1 . In general, the clustering decision problem is NP-complete and the clustering optimization problem is NP-hard, except in some simple cases involving either a constant number k of clusters or clustering in the 1-dimensional Euclidean space. In general, the results in Euclidean spaces are harder to derive than the results on graphs. When polynomial time algorithms exist, geometric problems tend to have faster solutions taking advantage of the geometric properties. However, none of the existing complexity theorems directly addresses the problem of clustering on the hypercube with respect to the Hamming distance.

To deal with the hypercube clustering problem one must first understand which quantities are allowed to go to infinity. If n is not allowed to go to infinity, then the number m of training examples is also bounded by 2^n and, since we are assuming $p < n$, there is no quantity that can scale. Thus by necessity we must have $n \rightarrow \infty$. We must also have $m \rightarrow \infty$. The case of interest for machine learning is when m is a low degree polynomial of n . Obviously the hypercube clustering problem is in NP, and it is a special case of clustering in \mathbb{R}^n . Thus the only important problem to be addressed is the reduction of a known NP-complete problem to a hypercube clustering problem.

For the reduction, it is natural to start from a known NP-complete graphical or geometric clustering problem. In both case, one must find ways to embed the original problem with its original metric into the hypercube with the Hamming distance. There are theorems for homeomorphic or squashed-embedding of graphs into the hypercube ([Hartman, 1976](#); [Winkler, 1983](#)), however these embeddings do not map the original dissimilarity function onto the the Hamming metric. Thus here we prefer to start from some of the known

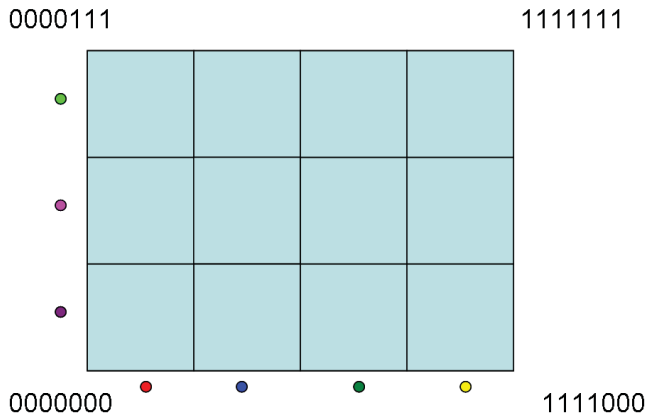


Figure 3: Embedding of a 3×4 Square Lattice onto \mathbb{H}^7 by Edge Coloring. All edges in the same row or column are given the same color. Each color corresponds to one of the dimensions of the 7-dimensional hypercube. For any pair of points, their Manhattan distance on the lattice is equal to the Hamming distance between their images in the 7-dimensional hypercube.

geometric results and use a strict cubical graph embedding. A graph is cubical if it is the subgraph of some hypercube \mathbb{H}^d for some d (Harary, 1988; Livingston and Stout, 1988). Although deciding whether a graph is cubical is NP-complete (Afrati et al., 1985), there is a theorem (Havel and Morávek, 1972) providing a necessary and sufficient condition for a graph to be cubical. A graph $G(V, E)$ is cubical and embeddable in \mathbb{H}^d if and only if it is possible to color the edges of G with d colors such that: (1) All edges incident with a common vertex are of different color; (2) In each path of G , there is some color that appears an odd number of times; and (3) In each cycle of G , no color appears an odd number of times. We can now state and prove the following theorem.

Theorem. Consider the following hypercube clustering problem:

Input: m binary vectors x_1, \dots, x_m of length n and an integer k .

Output: k binary vectors c_1, \dots, c_k of length n (the centroids) and a function f from $\{x_1, \dots, x_m\}$ to $\{c_1, \dots, c_k\}$ that minimizes the distortion $E = \sum_{t=1}^m \Delta(x_t, f(x_t))$ where Δ is the Hamming distance. The hypercube clustering problem is NP hard when $k \sim m^\epsilon$ ($\epsilon > 0$).

Proof. To sketch the reduction, we start from the problem of clustering m points in the plane \mathbb{R}^2 using cluster centroids and the L_1 distance, which is NP-complete (Megiddo and Supowit, 1984) by reduction from 3-SAT (Garey and Johnson, 1979) when $k \sim m^\epsilon$ ($\epsilon > 0$) (see, also related results in Mahajan et al. (2009) and Vattani (2010)). Without any loss of generality, we can assume that the points in these problems lie on the vertices of a square lattice. Using the theorem in Havel and Morávek (1972), one can show that a $n \times m$ square lattice in the plane can be embedded into \mathbb{H}^{n+m} . In fact, an explicit embedding is given in Figure 3. It is easy to check that the L_1 or Manhattan distance between any two points on the square lattice is equal to the corresponding Hamming distance in \mathbb{H}^{n+m} . This

polynomial reduction completes the proof that if the number of cluster satisfies $k = 2^p \sim m^\epsilon$, or equivalently $p \sim \epsilon \log_2 m \sim C \log n$, then the hypercube clustering problem associated with the Boolean autoencoder is NP-hard, and the corresponding decision problem NP-complete. If the numbers k of clusters is fixed and the centroids must belong to the training set, there are only $\binom{m}{k} \sim m^k$ possible choices for the centroids inducing the corresponding Voronoi clusters. This yields a trivial, albeit not efficient, polynomial time algorithm. When the centroids are not required to be in the training set, we conjecture also the existence of polynomial time algorithms by adapting the corresponding theorems in Euclidean space.

6. The Case $p \geq n$

When the hidden layer is larger than the input layer and $\mathbb{F} = \mathbb{G}$, there is an optimal 0-distortion solution involving the identity function. Thus this case is interesting only if additional constraints are added to the problem. These can come in the form of regularization, for instance to ensure sparsity of the hidden-layer representation, or restrictions on the classes of functions \mathcal{A} and \mathcal{B} , or noise in the hidden layer (see next section). When these constraints force the hidden layer to assume only k different values and $k < m$, for instance in the case of a sparse Boolean hidden layer, then the previous analyses hold and the problem reduces to a k clustering problem.

In this context of large hidden layers, in addition to vertical composition, there is also a natural horizontal composition for autoencoders that can be used to create large hidden layer representations (Figure 4) simply by horizontally combining autoencoders. Two (or more) autoencoders with architectures $n/p_1/n$ and $n/p_2/n$ can be trained and the hidden layers can be combined to yield an expanded hidden representation of size $p_1 + p_2$ that can then be fed to the subsequent layers of the overall architecture. Differences in the p_1 and p_2 hidden representations could be introduced by many different mechanisms, for instance using different learning algorithms, different initializations, different training samples, different learning rates, or different distortion measures. It is also possible to envision algorithms that incrementally add (or remove) hidden units to the hidden layer (Reed, 1993; Kwok and Yeung, 1997). In the linear case over \mathbb{R} , for instance, a first hidden unit can be trained to extract the first principal component, a second hidden unit can then be added to extract the second principal component, and so forth.

7. Other Generalizations

Within the general framework introduced here, other kinds of autoencoders can be considered. First, one can consider mixed autoencoders with different constraints on \mathbb{F} and \mathbb{G} , or different constraints on \mathcal{A} and \mathcal{B} . A simple example is when the input and output layers are real $\mathbb{F} = \mathbb{R}$ and the hidden layer is binary $\mathbb{G} = \{0, 1\}$ (and $\Delta = L_2^2$). It is easy to check that in this case, as long as $2^p = k < m$, the autoencoder aims at clustering the real data into k clusters and all the results obtained in the Boolean case are applicable with the proper adjustments. For instance, the centroid associated with a hidden state h should be the center of mass of the input vectors mapped onto h . In general, this mixed autoencoder is also NP hard when $k \sim m^\epsilon$ and, from a probabilistic view point, it corresponds to a mixture of k Gaussians model.

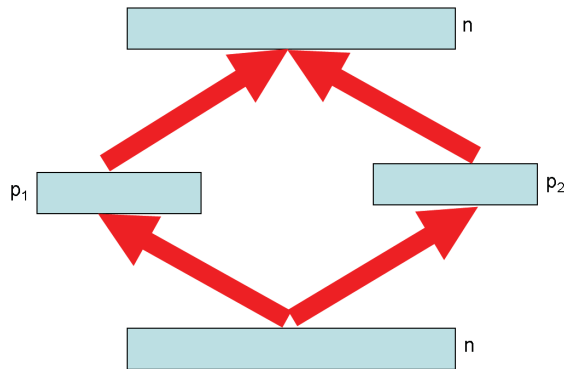


Figure 4: Horizontal Composition of Autoencoders to Expand the Hidden Layer Representation. The hidden layers of two separate autoencoders can be combined to yield a larger hidden representation of size $p_1 + p_2$ (see text).

A second natural direction is to consider autoencoders that are linear but over fields other than the real numbers, for instance over the field \mathbb{C} of complex numbers (Baldi et al., 2011), or over finite fields. For all these linear autoencoders, the Kernel of B plays an important role since inputs vectors are basically clustered modulo this kernel. These autoencoders are not without theoretical and practical interests. Consider the linear autoencoder over the Galois field with two elements $GF(2) = \mathbb{F}_2$. It is easy to see that this is a special case of the Boolean autoencoder, where the Boolean functions are restricted to parity functions. This autoencoder can also be seen as implementing a linear code (McEliece, 1977). When there is noise in the ‘transmission’ of the hidden layer and $p > n$, one can consider solutions where n units in the hidden layer correspond to the identity function and the remaining $p - n$ units implement additional parity check bits that are linearly computed from the input and used for error correction. Thus all well known linear codes, such as Hamming or Reed-Solomon codes, can be viewed within this linear autoencoder framework. While the linear autoencoder over \mathbb{F}_2 will be discussed elsewhere, it is worth noting that it is likely to yield an NP-hard problem, as for the case of the unrestricted Boolean autoencoder. This can be seen by considering that finding the minimum (non-zero) weight vector in the kernel of a binary matrix, or the radius of a code, are NP-complete problems (McEliece and van Tilborg, 1978; Frances and Litman, 1997). A simple overall classification of autoencoders is given in Figure 5.

8. Discussion

Studying the linear and Boolean autoencoder in detail enables one to gain a general perspective on autoencoders, define key properties that are shared by different autoencoders and that ought to be checked systematically in any new kind of autoencoder (e.g. group invariances, clustering, recycling stability). The general perspective also shows the intimate connections between autoencoders and information and coding theory: (1) autoencoders

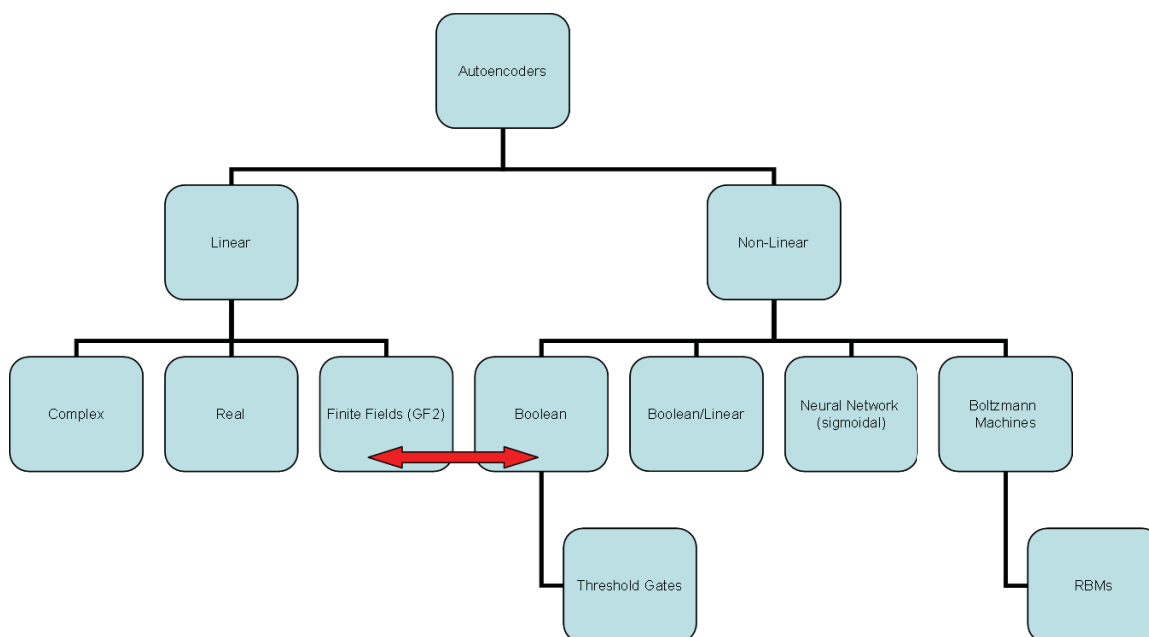


Figure 5: Simple Autoencoder Classification.

with $n < p$ and a noisy hidden layer correspond to the classical noisy channel transmission and coding problem, with linear code over finite fields as a special case; (2) autoencoders with $n > p$ correspond to compression which can be lossy when the number of states in the hidden layer is less than the number of training examples ($m > k$) and lossless otherwise ($m < k$).

When $n > p$ the general emerging picture is that autoencoders learning is in general NP-complete¹ except in simple but important cases (e.g. linear over \mathbb{R} , Boolean with fixed k) and that in essence all autoencoders are performing some form of clustering suggesting a unified view of different forms of unsupervised learning, where Hebbian learning, autoencoders, and clustering are three faces of the same die. While autoencoders and Hebbian rules provide unsupervised learning implementations, it is clustering that provides the basic conceptual operation that underlies them. In addition, it is important to note that in general clustering objects without providing a label for each cluster is useless for subsequent higher levels of processing. In addition to clustering, autoencoders provide a label for each cluster through the activity in the hidden layer and thus elegantly address both the clustering and labeling problems at the same time.

RBMs and their efficient contrastive learning algorithm may provide an efficient form of autoencoder and autoencoder learning, but it is doubtful that there is anything special about RBMs at a deeper conceptual level. Thus it ought to be possible to derive results comparable to those described in Hinton et al. (2006) and Hinton and Salakhutdinov (2006) by stacking other kinds of autoencoders, and more generally by hierarchically stacking a series of clustering algorithms using vertical composition, perhaps also in combination with

1. RBM learning is NP-complete by similarity with minimizing a quadratic form over the hypercube.

horizontal composition. As pointed out in the previous sections, it is easy to add a top layer for supervised regression or classification tasks on top of the hierarchical clustering stack. In aggregate, these results suggest that: (1) the so-called deep architectures may in fact have a non-trivial but constant (or logarithmic) depth, which is also consistent with what is observed in sensory neuronal circuits; (2) the fundamental unsupervised operation behind deep architectures, in one form or the other, is clustering, which is composable both horizontally and vertically—in this view, clustering may emerge as the central and almost sufficient ingredient for building intelligent systems; and (3) the generalization properties of deep architectures may be easier to understand when ignoring many of the hardware details, in terms of the most simple forms of autoencoders (e.g Boolean), or in terms of the more fundamental underlying clustering operations.

Acknowledgments

Work in part supported by grants NSF IIS-0513376, NIH LM010235, and NIH-NLM T15 LM07443 to PB.

References

- F. Afrati, C.H. Papadimitriou, and G. Papageorgiou. The complexity of cubical graphs. *Information and control*, 66(1-2):53–60, 1985.
- P. Baldi and K. Hornik. Neural networks and principal component analysis: Learning from examples without local minima. *Neural Networks*, 2(1):53–58, 1988.
- P. Baldi, S. Forouzan, and Z. Lu. Complex-Valued Autoencoders. *Neural Networks*, 2011. Submitted.
- Yoshua Bengio and Yann LeCun. Scaling learning algorithms towards AI. In L. Bottou, O. Chapelle, D. DeCoste, and J. Weston, editors, *Large-Scale Kernel Machines*. MIT Press, 2007.
- H. Boursard and Y. Kamp. Auto-association by multilayer perceptrons and singular value decomposition. *Biological cybernetics*, 59(4):291–294, 1988. ISSN 0340-1200.
- P. Clote and E. Kranakis. *Boolean functions and computation models*. Springer Verlag, 2002.
- Dumitru Erhan, Yoshua Bengio, Aaron Courville, Pierre-Antoine Manzagol, Pascal Vincent, and Samy Bengio. Why does unsupervised pre-training help deep learning? *Journal of Machine Learning Research*, 11:625–660, February 2010.
- M. Frances and A. Litman. On covering problems of codes. *Theory of Computing Systems*, 30(2):113–119, 1997.
- B.J. Frey and D. Dueck. Clustering by passing messages between data points. *Science*, 315(5814):972, 2007.
- M.R. Garey and D.S. Johnson. *Computers and Intractability*. Freeman San Francisco, 1979.

- F. Harary. Cubical graphs and cubical dimensions. *Computers & Mathematics with Applications*, 15(4):271–275, 1988.
- J. Hartman. The homeomorphic embedding of K_n in the m -cube* 1. *Discrete Mathematics*, 16(2):157–160, 1976.
- I. Havel and J. Morávek. B -valuations of graphs. *Czechoslovak Mathematical Journal*, 22(2):338–351, 1972.
- D.O. Hebb. The organization of behavior: A neurophysiological study. *Wiley-Interscience, New York*, 1949.
- G.E. Hinton and R.R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504, 2006.
- G.E. Hinton, S. Osindero, and Y.W. Teh. A fast learning algorithm for deep belief nets. *Neural Computation*, 18(7):1527–1554, 2006.
- T. Kwok and D. Yeung. Constructive Algorithms for Structure Learning in Feedforward Neural Networks for Regression Problems. *IEEE Transactions on Neural Networks*, 8: 630–645, 1997.
- M. Livingston and Q.F. Stout. Embeddings in hypercubes. *Mathematical and Computer Modelling*, 11:222–227, 1988.
- M. Mahajan, P. Nimbhorkar, and K. Varadarajan. The planar k -means problem is NP-hard. *WALCOM: Algorithms and Computation*, pages 274–285, 2009.
- R. McEliece and H. van Tilborg. On the inherent intractability of certain coding problems(Corresp.). *IEEE Transactions on Information Theory*, 24(3):384–386, 1978.
- R. J. McEliece. *The Theory of Information and Coding*. Addison-Wesley Publishing Company, Reading, MA, 1977.
- N. Megiddo and K.J. Supowit. On the complexity of some common geometric location problems. *SIAM J. COMPUT.*, 13(1):182–196, 1984.
- G. Montufar and N. Ay. Refinements of Universal Approximation Results for Deep Belief Networks and Restricted Boltzmann Machines. *Neural Computation*, pages 1–14, 2011. ISSN 0899-7667.
- E. Oja. Simplified neuron model as a principal component analyzer. *Journal of mathematical biology*, 15(3):267–273, 1982.
- R. Reed. Pruning algorithms-a survey. *Neural Networks, IEEE Transactions on*, 4(5): 740–747, 1993.
- D.E. Rumelhart, G.E. Hinton, and R.J. Williams. Learning internal representations by error propagation. In *Parallel Distributed Processing. Vol 1: Foundations*. MIT Press, Cambridge, MA, 1986.

- JL Slagle, CL Chang, and SR Heller. A clustering and data reorganization algorithm. *IEEE Transactions on Systems, Man and Cybernetics*, 5:121–128, 1975.
- I. Sutskever and G.E. Hinton. Deep, narrow sigmoid belief networks are universal approximators. *Neural Computation*, 20(11):2629–2636, 2008.
- A. Vattani. A simpler proof of the hardness of k-means clustering in the plane. *UCSD Technical Report*, 2010.
- P.M. Winkler. Proof of the squashed cube conjecture. *Combinatorica*, 3(1):135–139, 1983.