

# A Review of Unsupervised Feature Learning and Deep Learning for Time-Series Modeling

Martin Längkvist<sup>a,\*</sup>, Lars Karlsson<sup>a</sup>, Amy Loutfi<sup>a</sup>

<sup>a</sup>*Applied Autonomous Sensor Systems, School of Science and Technology, Örebro University, SE-701 82, Örebro, Sweden*

---

## Abstract

This paper gives a review of the recent developments in deep learning and unsupervised feature learning for time-series problems. While these techniques have shown promise for modeling static data, such as computer vision, applying them to time-series data is gaining increasing attention. This paper overviews the particular challenges present in time-series data and provides a review of the works that have either applied time-series data to unsupervised feature learning algorithms or alternatively have contributed to modifications of feature learning algorithms to take into account the challenges present in time-series data.

*Keywords:* time-series, unsupervised feature learning, deep learning

---

## 1. Introduction and Background

Time is a natural element that is always present when the human brain is learning tasks like language, vision and motion. Most real-world data has a temporal component, whether it is measurements of natural processes

---

\*Corresponding author

*Email addresses:* martin.langkvist@oru.se (Martin Längkvist), lars.karlsson@oru.se (Lars Karlsson), amy.loutfi@oru.se (Amy Loutfi)

5 (weather, sound waves) or man-made (stock market, robotics). Analysis of  
6 time-series data has been the subject of active research for decades (Keogh  
7 and Kasetty, 2002; Dietterich, 2002) and is considered by Yang and Wu  
8 (2006) as one of the top 10 challenging problems in data mining due to  
9 its unique properties. Traditional approaches for modeling sequential data  
10 include the estimation of parameters from an assumed time-series model,  
11 such as autoregressive models (Lütkepohl, 2005) and Linear Dynamical Sys-  
12 tems (LDS) (Luenberger, 1979), and the popular Hidden Markov Model  
13 (HMM) (Rabiner and Juang, 1986). The estimated parameters can then  
14 be used as features in a classifier to perform classification. However, more  
15 complex, high-dimensional, and noisy real-world time-series data cannot be  
16 described with analytical equations with parameters to solve since the dy-  
17 namics are either too complex or unknown (Taylor, 2009) and traditional  
18 shallow methods, which contain only a small number of non-linear opera-  
19 tions, do not have the capacity to accurately model such complex data.

20 In order to better model complex real-world data, one approach is to  
21 develop robust features that capture the relevant information. However, de-  
22 veloping domain-specific features for each task is expensive, time-consuming,  
23 and requires expertise of the data. The alternative is to use unsupervised  
24 feature learning (Bengio and LeCun, 2007; Bengio et al., 2012; Erhan et al.,  
25 2010) in order to learn a layer of feature representations from unlabeled data.  
26 This has the advantage that the unlabeled data, which is plentiful and easy  
27 to obtain, is utilized and that the features are learned from the data instead  
28 of being hand-crafted. Another benefit is that these layers of feature repre-  
29 sentations can be stacked to create deep networks, which are more capable

30 of modeling complex structures in the data. Deep networks have been used  
31 to achieve state-of-the-art results on a number of benchmark data sets and  
32 for solving difficult AI tasks. However, much focus in the feature learning  
33 community has been on developing models for static data and not so much  
34 on time-series data.

35 In this paper we review the variety of feature learning algorithms that  
36 has been developed to explicitly capture temporal relationships as well as the  
37 various time-series problems that they have been used on. The properties of  
38 time-series data will be discussed in Section 2 followed by an introduction to  
39 unsupervised feature learning and deep learning in Section 3. An overview  
40 of some common time-series problems and previous work using deep learning  
41 is given in Section 4. Finally, conclusions are given in Section 5.

## 42 **2. Properties of time-series data**

43 Time-series data consists of sampled data points taken from a continuous,  
44 real-valued process over time. There are a number of characteristics of time-  
45 series data that make it different from other types of data.

46 Firstly, the sampled time-series data often contain much noise and have  
47 high dimensionality. To deal with this, signal processing techniques such  
48 as dimensionality reduction techniques, wavelet analysis or filtering can be  
49 applied to remove some of the noise and reduce the dimensionality. The use  
50 of feature extraction has a number of advantages (Nanopoulos et al., 2001).  
51 However, valuable information could be lost and the choice of features and  
52 signal processing techniques may require expertise of the data.

53 The second characteristics of time-series data is that it is not certain

54 that there are enough information available to understand the process. For  
55 example, in electronic nose data, where an array of sensors with various  
56 selectivity for a number of gases are combined to identify a particular smell,  
57 there is no guarantee that the selection of sensors actually are able to identify  
58 the target odour. In financial data when observing a single stock, which only  
59 measures a small aspect of a complex system, there is most likely not enough  
60 information in order to predict the future (Fama, 1965).

61 Further, time-series have an explicit dependency on the time variable.  
62 Given an input  $x(t)$  at time  $t$ , the model predicts  $y(t)$ , but an identical input  
63 at a later time could be associated with a different prediction. To solve this  
64 problem, the model either has to include more data input from the past or  
65 must have a memory of past inputs. For long-term dependencies the first ap-  
66 proach could make the input size too large for the model to handle. Another  
67 challenge is that the length of the time-dependencies could be unknown.

68 Many time-series are also non-stationary, meaning that the characteristics  
69 of the data, such as mean, variance, and frequency, changes over time. For  
70 some time-series data, the change in frequency is so relevant to the task that it  
71 is more beneficial to work in the frequency-domain than in the time-domain.

72 Finally, there is a difference between time-series data and other types of  
73 data when it comes to invariance. In other domains, for example computer  
74 vision, it is important to have features that are invariant to translations,  
75 rotations, and scale. Most features used for time-series need to be invariant  
76 to translations in time.

77 In conclusion, time-series data is high-dimensional and complex with  
78 unique properties that make them challenging to analyze and model. There

79 is a large interest in representing the time-series data in order to reduce the  
80 dimensionality and extract relevant information. The key for any success-  
81 ful application lies in choosing the right representation. Various time-series  
82 problems contain different degrees of the properties discussed in this section  
83 and prior knowledge or assumptions about these properties is often infused  
84 in the chosen model or feature representation. There is an increasing in-  
85 terest in learning the representation from unlabeled data instead of using  
86 hand-designed features. Unsupervised feature learning have shown to be  
87 successful at learning layers of feature representations for static data sets  
88 and can be combined with deep networks to create more powerful learning  
89 models. However, the feature learning for time-series data have to be mod-  
90 ified in order to adjust for the characteristics of time-series data in order to  
91 capture the temporal information as well.

### 92 **3. Unsupervised feature learning and deep learning**

93 This section presents both models that are used for unsupervised feature  
94 learning and models and techniques that are used for modeling temporal  
95 relations. The advantage of learning features from unlabeled data is that the  
96 plentiful unlabeled data can be utilized and that potentially better features  
97 than hand-crafted features can be learned. Both these advantages reduce the  
98 need for expertise of the data.

#### 99 *3.1. Restricted Boltzmann Machine*

100 The Restricted Boltzmann Machines (RBM) (Hinton et al., 2006; Hinton  
101 and Salakhutdinov, 2006; Lee et al., 2008) is a generative probabilistic model  
102 between input units (visible),  $\mathbf{x}$ , and latent units (hidden),  $\mathbf{h}$ , see Figure 1.

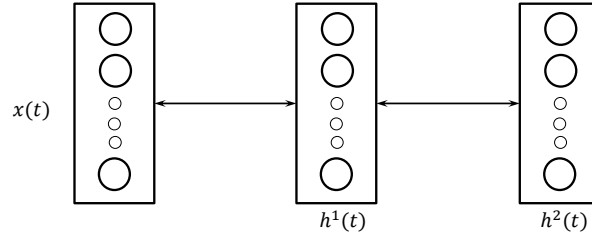


Figure 1: A 2-layer RBM for static data. The visible units  $x$  are fully connected to the first hidden layer  $h^1$ .

103 The visible and hidden units are connected with a weight matrix,  $\mathbf{W}$  and  
 104 have bias vectors  $\mathbf{c}$  and  $\mathbf{b}$ , respectively. There are no connections among  
 105 the visible and hidden units. The RBM can be used to model static data.  
 106 The energy function and the joint distribution for a given visible and hidden  
 107 vector is defined as:

$$E(\mathbf{x}, \mathbf{h}) = \mathbf{h}^T \mathbf{W} \mathbf{x} + \mathbf{b}^T \mathbf{h} + \mathbf{c}^T \mathbf{x} \quad (1)$$

$$P(\mathbf{x}, \mathbf{h}) = \frac{1}{Z} \exp^{E(\mathbf{x}, \mathbf{h})} \quad (2)$$

108 where  $Z$  is the partition function that ensures that the distribution is nor-  
 109 malized. For binary visible and hidden units, the probability that hidden  
 110 unit  $h_j$  is activated given visible vector  $x$  and the probability that visible  
 111 unit  $x_i$  is activated given hidden vector  $h$  are given by:

$$P(h_j|\mathbf{x}) = \sigma(b_j + \sum_i W_{ij}x_i) \quad (3)$$

$$P(x_i|\mathbf{h}) = \sigma(c_i + \sum_j W_{ij}h_j) \quad (4)$$

112 where  $\sigma(\cdot)$  is the activation function. The logistic function,  $\sigma(x) = \frac{1}{1+e^{-x}}$ ,  
 113 is a common choice for the activation function. The parameters  $W$ ,  $b$ , and  
 114  $v$ , are trained to minimize the reconstruction error using contrastive diver-  
 115 gence (Hinton, 2002). The learning rule for the RBM is:

$$\frac{\partial \log P(\mathbf{x})}{\partial W_{ij}} \approx \langle x_i h_j \rangle_{data} - \langle x_i h_j \rangle_{recon} \quad (5)$$

116 where  $\langle \cdot \rangle$  is the average value over all training samples. Several RBMs can  
 117 be stacked to produce a deep belief network (DBN). In a deep network, the  
 118 activation of the hidden units in the first layer is the input to the second  
 119 layer.

### 120 3.2. Conditional RBM

121 An extension of RBM that models multivariate time-series data is the  
 122 conditional RBM (cRBM), see Figure 2. A similar model is the Temporal  
 123 RBM (Sutskever and Hinton, 2006). The cRBM consists of auto-regressive  
 124 weights that model short-term temporal structures, and connections between

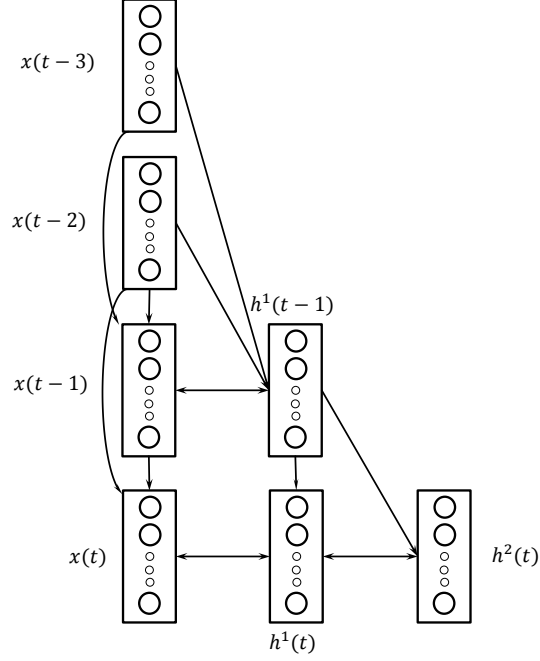


Figure 2: A 2-layer conditional RBM for time-series data. The model order for the first and second layer is 3 and 2, respectively.

125 past visible units to the current hidden units. The bias vectors in a cRBM  
 126 depend on previous visible units and are defined as:

$$b_j^* = b_j + \sum_{i=1}^n B_i x(t-i) \quad (6)$$

$$c_i^* = c_i + \sum_{j=1}^n A_j x(t-j) \quad (7)$$

127 where  $A_i$  is the auto-regressive connections between visible units at time  $t-i$   
 128 and current visible units,  $B_i$  is the weight matrix connecting visible layer at



time  $t - i$  to the current hidden units. The model order is defined by the constant  $n$ . The probabilities for going up or down a layer are:

$$P(h_j|\mathbf{x}) = \sigma \left( b_j + \sum_i W_{ij}x_i + \sum_k \sum_i B_{ijk}x_i(t - k) \right) \quad (8)$$

$$P(x_i|\mathbf{h}) = \sigma \left( c_i + \sum_j W_{ij}h_j + \sum_k \sum_i A_{ijk}x_i(t - k) \right) \quad (9)$$

The parameters  $\theta = \{W, b, c, A, B\}$ , are trained using contrastive divergence. Just like a RBM, the cRBM can also be used as a module to create deep networks.

### 3.3. Gated RBM

The Gated Restricted Boltzmann Machine (GRBM) (Memisevic and Hinton, 2007) is another extension of the RBM that models the transition between two input vectors. The GRBM models a weight tensor,  $W_{ijk}$ , between the input,  $\mathbf{x}$ , the output,  $\mathbf{y}$ , and latent variables,  $\mathbf{z}$ . The energy function is defined as:

$$E(\mathbf{y}, \mathbf{z}; \mathbf{x}) = - \sum_{ijk} W_{ijk}x_i y_j z_k - \sum_k b_k z_k - \sum_j c_j y_j \quad (10)$$

where  $\mathbf{b}$  and  $\mathbf{c}$  are the bias vectors for  $\mathbf{x}$  and  $\mathbf{y}$ , respectively. The conditional probability of the transformation and the output image given the input image is:

$$p(\mathbf{y}, \mathbf{z}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp(-E(\mathbf{y}, \mathbf{z}; \mathbf{x})) \quad (11)$$

143 where  $Z(\mathbf{x})$  is the partition function. Luckily, this quantity does not need to  
 144 be computed to perform inference or learning. The probability that hidden  
 145 unit  $z_i$  is activated given  $\mathbf{x}$  and  $\mathbf{y}$  is given by:

$$P(z_k = 1|\mathbf{x}, \mathbf{y}) = \sigma\left(\sum_{ij} W_{ijk}x_iy_j + b_k\right) \quad (12)$$

146 Learning the parameters is performed with an approximation method of the  
 147 gradient called contrastive divergence (Hinton, 2002). Each latent variable  
 148  $z_k$  learns a simple transformation that together are combined the represent  
 149 the full transformation. By fixating a learned transformation  $\mathbf{z}$  and given  
 150 an input image  $\mathbf{x}$ , the output image  $\mathbf{y}$  is the selected transformation applied  
 151 to the input image. Similarly, for a fixed input image  $\mathbf{x}$ , a given image  $\mathbf{y}$   
 152 creates a RBM that learns the transformation  $\mathbf{z}$  by reconstructing  $\mathbf{y}$ . These  
 153 properties could not be achieved with a regular RBM with input units sim-  
 154 ply being the concatenated images  $\mathbf{x}$  and  $\mathbf{y}$  since the latent variables would  
 155 only learn the spatial information for that particular image pair and not the  
 156 general transformation. The large number of parameters due to the weight  
 157 tensor makes it impractical for large image sizes. A factored form of the  
 158 three-way tensor has been proposed to reduce the number of parameters to  
 159 learn (Memisevic and Hinton, 2010).

### 160 3.4. *Auto-encoder*

161 A model that does not have a partition function is the auto-encoder (Ran-  
 162 zato et al., 2006; Bengio et al., 2007; Bengio, 2007), see Figure 3. The auto-  
 163 encoder was first introduced as a dimensionality reduction algorithm. In  
 164 fact, a basic linear auto-encoder learns essentially the same representation as

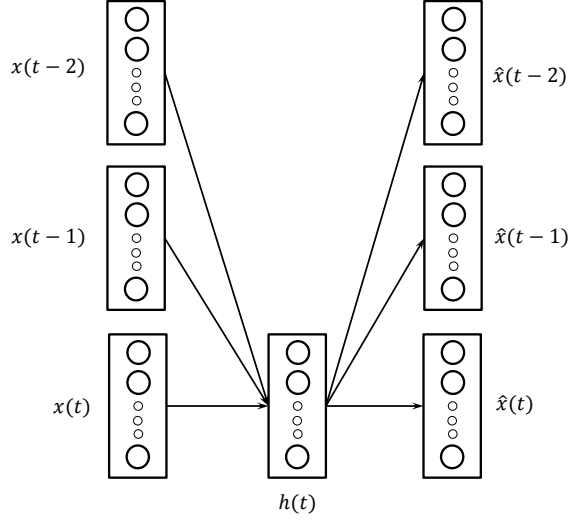


Figure 3: A 1-layer auto-encoder for static time-series input. The input is the concatenation of current and past frames of visible data  $x$ . The reconstruction of  $x$  is denoted  $\hat{x}$ .

165 a Principal Component Analysis (PCA). The layers of visible units,  $\mathbf{x}$ , hid-  
 166 den units,  $\mathbf{h}$ , and the reconstruction of the visible units,  $\hat{\mathbf{x}}$ , are connected via  
 167 weight matrices  $\mathbf{W}^1$  and  $\mathbf{W}^2$  and the hidden layer and reconstruction layer  
 168 have bias vectors  $\mathbf{b}^1$  and  $\mathbf{b}^2$ , respectively. It is common in auto-encoders  
 169 to have tied weights, that is,  $\mathbf{W}^2 = (\mathbf{W}^1)^T$ . This works as a regularizer  
 170 as it constrains the allowed parameter space and reduces the number of pa-  
 171 rameters to learn (Bengio et al., 2012). The feed-forward activations are  
 172 calculated as:

$$h_j = \sigma(\sum_i W_{ji}^1 x_i + b_j^1) \quad (13)$$

$$\hat{x}_i = \sigma(\sum_j W_{ij}^2 h_j + b_i^2) \quad (14)$$

where  $\sigma(\cdot)$  is the activation function. As with the RBM, a common choice is the logistic activation function. The cost function to be minimized is expressed as:

$$J(\theta) = \frac{1}{2N} \sum_n \sum_i (x_i^{(n)} - \hat{x}_i^{(n)})^2 + \frac{\lambda}{2} \sum_l \sum_i \sum_j (W_{ij}^l)^2 + \beta \sum_l \sum_j KL(\rho || p_j^l) \quad (15)$$

where  $p_j^l$  is the mean activation for unit  $j$  in layer  $l$ ,  $\rho$  is the desired mean activation, and  $N$  is the number of training examples. KL is the Kullback-Leibler (KL) divergence which is defined as  $KL(\rho || p_j^l) = \rho \log \frac{\rho}{p_j^l} + (1 - \rho) \log \frac{1-\rho}{1-p_j^l}$ . The first term is the square root error term that will minimize the reconstruction error. The second term is the L2 weight decay term that will keep the weight matrices close to zero. Finally, the third term is the sparsity penalty term and encourages each unit to only be partially activated as specified by the hyperparameter  $\rho$ . The inclusion of these regularization terms prevents the trivial learning of a 1-to-1 mapping of the input to the hidden units. A difference between auto-encoders and RBMs is that RBMs do not require such regularization because the use of stochastic binary hidden units acts as a very strong regularizer (Hinton, 2012). However, it is not uncommon to introduce an extra sparsity constraint for RBMs (Lee et al., 2008).

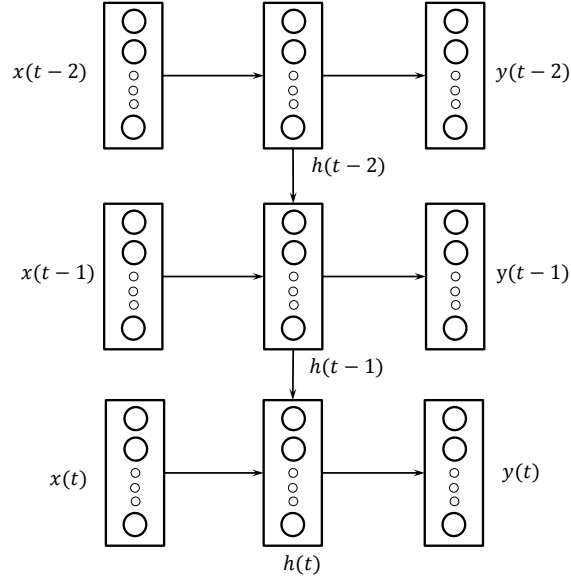


Figure 4: A Recurrent Neural Network (RNN). The input  $x$  is transformed to the output representation  $y$  via the hidden units  $h$ . The hidden units have connections from the input values of the current time frame and the hidden units from the previous time frame.

191 A model that have been used for modeling sequential data is the Recur-  
 192 rent Neural Network (RNN) (Hüsken and Stagge, 2003). Generally, an RNN  
 193 is obtained from the feedforward network by connecting the neurons' output  
 194 to their inputs, see Figure 4. The short-term time-dependency is modeled by  
 195 the hidden-to-hidden connections without using any time delay-taps. They  
 196 are usually trained iteratively via a procedure known as backpropagation-  
 197 through-time (BPTT). RNNs can be seen as very deep networks with shared  
 198 parameters at each layer when unfolded in time. This results in the prob-

199 lem of vanishing gradients (Pascanu et al., 2012) and has motivated the  
200 exploration of second-order methods for deep architectures (Martens and  
201 Sutskever, 2012) and unsupervised pre-training. An overview of strategies  
202 for training RNNs is provided by Sutskever (2012). A popular extension is  
203 the use of the purpose-built Long-short term memory cell (Hochreiter and  
204 Schmidhuber, 1997) that better finds long-term dependencies.

### 205 3.6. *Deep learning*

206 The models presented in this section use a non-linear activation function  
207 on the hidden units. This non-linearity enables a more expressive model that  
208 can learn more abstract representations when multiple modules are stacked  
209 on top of each other to form a deep network (if linear features would be  
210 stacked the result would still be a linear operation). The goal of a deep net-  
211 work is to build features at the lower layers that will disentangle the factors  
212 of variations in the input data and then combine these representations at  
213 the higher layers. It has been proposed that a deep network will generalize  
214 better because it has a more compact representation (Le Roux and Bengio,  
215 2008). However, the difficulty with training multiple layers of hidden units  
216 lies in the problem of vanishing gradients when the error signal is backpropa-  
217 gated (Bengio et al., 1994). This can be solved by doing unsupervised greedy  
218 layer-wise pre-training of each layer. This acts as an unusual form of regular-  
219 ization (Erhan et al., 2010) that avoids poor local minima and gives a better  
220 initialization than a random initialization (Bengio et al., 2012). However,  
221 the importance of parameter initialization is not as crucial as other factors  
222 such as input connections and architecture (Saxe et al., 2011).

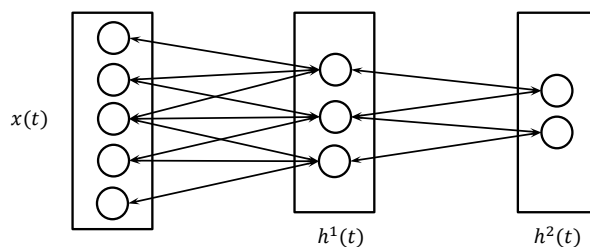


Figure 5: A 2-layer convolutional neural network.

224 A technique that is particularly interesting for high-dimensional data,  
 225 such as images and time-series data, is convolution. In a convolutional set-  
 226 ting, the hidden units are not fully connected to the input but instead di-  
 227 vided into locally connected segments, see Figure 5. Convolution has been  
 228 applied to both RBMs and auto-encoders to create convolutional RBMs (con-  
 229 vRBM) (Lee et al., 2009b,a) and convolutional auto-encoders (convAE) (Masci  
 230 et al., 2011). A Time-Delay Neural Network (TDNN) is a specialization of  
 231 Artificial Neural Networks (ANN) that exploits the time structure of the  
 232 input by performing convolutions on overlapping windows.

233 A common operator used together with convolution is pooling, which  
 234 combines nearby values in input or feature space through a max, average or  
 235 histogram operator. The purpose of pooling is to achieve invariance to small  
 236 local distortions and reduce the dimensionality of the feature space. The work  
 237 by Lee et al. (2009a) introduces probabilistic max-pooling in the context  
 238 of convolutional RBMs. The Space-Time DBN (ST-DBN) (Bo Chen and  
 239 de Freitas, 2010) uses convolutional RBMs together with a spatial pooling  
 240 layer and a temporal pooling layer to build invariant features from spatio-  
 241 temporal data.

### 242 3.8. Temporal coherence

243 There are a number of other ways besides the architectural structure  
 244 that can be used to capture temporal coherence in data. One way is to in-  
 245 troduce a smoothness penalty on the hidden variables in the regularization.  
 246 This is done by minimizing the changes in the hidden unit activations from  
 247 one frame to the next by  $\min |h(t) - h(t - 1)|$ . The motivation behind this  
 248 is that for sequential data the hidden unit activations should not change  
 249 much if the time-dependent data is fed to the model in a chronological or-  
 250 der. Other strategies include penalizing the squared difference, slow feature  
 251 analysis (Wiskott and Sejnowski, 2002), or as a function of other factors, for  
 252 example the change in the input data in order to adapt to both slow and  
 253 rapid changing input data.

254 Temporal coherence is related to invariant feature representations since  
 255 both methods want to achieve small changes in the feature representation for  
 256 small changes in the input data. It is suggested in Hinton et al. (2011) that  
 257 the pose parameters and affine transformations should be modeled instead



258 of using invariant feature representations. In that case, temporal coherence  
259 should be over a group of numbers, such as the position and pose of the  
260 object rather than a single scalar. This could for example be achieved using  
261 a structured sparsity penalty (Kavukcuoglu et al., 2009).

### 262 3.9. Hidden Markov Model

263 The Hidden Markov Model (HMM) (Rabiner and Juang, 1986) is a pop-  
264 ular model for modeling sequential data and is defined by two probability  
265 distributions. The first one is the transition distribution  $P(y_t|y_{t-1})$ , which  
266 defines the probability of going from one hidden state  $y$  to the next hidden  
267 state. The second one is the observation distribution  $P(x_t|y_t)$ , which defines  
268 the relation between observed  $x$  values and hidden  $y$  states. One assumption  
269 is that these distributions are stationary. However, the main problem with  
270 HMMs are that they require a discrete state space, often have unrealistic  
271 independence assumptions, and have a limited representational capacity of  
272 their hidden states (Mohamed and Hinton, 2010). HMMs require  $2^N$  hidden  
273 states in order to model  $N$  bits of information about the past history.

### 274 3.10. Summary

275 Table 1 gives a summary of the briefly presented models in this section.  
276 The first column indicates whether the model is capable of capturing tem-  
277 poral relations. A model that captures temporal relations does so by having  
278 a memory of past inputs. The memory of a model, indicated in the second  
279 column, means how many steps back in time an input have on the current  
280 frame. Without the temporal order, any permutation of the feature sequence

281 would yield the same distribution (Humphrey et al., 2013). The implemen-  
 282 tation of a memory is performed differently between the models. In a cRBM,  
 283 delay taps are used to create a short-term dependency on past visible units.  
 284 The long-term dependency comes from modeling subsequent layers. This  
 285 means that the length of the memory for a cRBM is increased for each added  
 286 layer. The model order for a cRBM in one layer is typically below 5 for input  
 287 sizes around 50. A decrease in the input size would allow a higher model  
 288 order. In an RNN, hidden units in the current time frame are affected by the  
 289 state of the hidden units in the previous time frame. This can create a ripple  
 290 effect with a duration of potentially infinite time frames. On the other hand,  
 291 this ripple effect can be prevented by using a forget gate (Gers et al., 2000).  
 292 The use of Long-short term memory (Hochreiter and Schmidhuber, 1997) or  
 293 hessian-free optimizer (Martens and Sutskever, 2012) can produce recurrent  
 294 networks that has a memory of over 100 time steps. The Gated RBM and  
 295 the convolutional GRBM models transitions between pairs of input vectors  
 296 so the memory for these models is 2. The Space-Time DBN (Bo Chen and  
 297 de Freitas, 2010) models 6 sequences of outputs from the spatial pooling  
 298 layer, which is a longer memory than GRBM, but using a lower input size.

299 The last column in Table 1 indicates if the model is generative (as op-  
 300 posed to discriminative). A generative model can generate observable data  
 301 given a hidden representation and this ability is mostly used for generating  
 302 synthetic data of future time steps. Even though the auto-encoder is not  
 303 generative, a probabilistic interpretation can be made using auto-encoder  
 304 scoring (Kamyshanska and Memisevic, 2013; Bengio et al., 2013).

305 For selecting a model for a particular problem, a number of questions

Table 1: A summary of commonly used models for feature learning.

Method	Temporal relation	Memory	Typical input size	Generative
RBM	-	-	10-1000	✓
AE	-	-	10-1000	-
RNN	✓	1-100	50-1000	✓
cRBM	✓	2-5	50	✓
TDNN	✓	2-5	5-50	-
ANN	-	-	10-1000	-
GRBM	✓	2	<64x64	✓
ConvGRBM	✓	2	>64x64	✓
ConvRBM	-	-	>64x64	✓
ConvAE	-	-	>64x64	-
ST-DBN	✓	2-6	10x10	✓

306 should be taken into consideration: (1) Use a generative or discriminative  
 307 model? (2) What are the properties of the data? and (3) How large is the  
 308 input size? A generative model is preferred if the trained model should be  
 309 used for synthesizing new data or prediction tasks where partial input data  
 310 (data at  $t + 1$ ) need to be reconstructed. If the task is to do classification,  
 311 a discriminative model is sufficient. A discriminative model will attempt to  
 312 model the training data even if that data is noisy while a generative model  
 313 will simply assign a low probability for outliers. This makes a generative  
 314 model more robust for noisy inputs and a better outlier detector. There is  
 315 also the factor of training time. Generative models use Gibbs sampling to  
 316 approximate the derivatives for each parameter update while a discrimina-  
 317 tive model calculates the exact gradients in one iteration. However, if the  
 318 simulation time is an issue, it is a good idea to look for hardware solutions

319 or the choice of optimization method before considering which method is the  
320 fastest. When the combination of input size, model parameters, and number  
321 of training examples in one training batch is large, the training time could  
322 be decreased by performing the parameter updates on a GPU instead of the  
323 CPU. For large-scale problems, i.e., the number of training examples is large,  
324 it is recommended to use stochastic gradient descent instead of L-BFGS or  
325 conjugate gradient descent as optimization method (Bottou, 2010). Further-  
326 more, if the data has a temporal structure it is not recommended to treat  
327 the input data as a feature vector since this will discard the temporal in-  
328 formation. Instead, a model that inherently models temporal relations or  
329 incorporates temporal coherence (by regularization or temporal pooling) in  
330 a static model is a better approach. For high-dimensional problems, like  
331 images which have a pictorial structure, it may be appropriate to use convo-  
332 lution. The use of pooling further decreases the number of dimensions and  
333 introduces invariance for small translations of the input data.

#### 334 **4. Classical time-series problems**

335 In this section we will highlight some common time-series problems and  
336 the models that have been used to address them in the literature. We will fo-  
337 cus on complex problems that require the use of models with hidden variables  
338 for feature representation and where the representations are fully or partially  
339 learned from unlabeled data. A summary of the classical time-series problems  
340 that will be presented in this section is given in Table 2.



Figure 6: Four images from the KTH action recognition data set of a person running at frame 100, 105, 110, and 115. The KTH data set also contains videos of walking, jogging, boxing, hand waving, and handclapping.

#### 341 4.1. Videos

342 Video data are series of images over time (spatio-temporal data) and can  
 343 therefore be viewed as high-dimensional time-series data. Figure 6 shows  
 344 a sequence of images from the KTH activity recognition data set<sup>1</sup>. The  
 345 traditional approach to modeling video streams is to treat each individual  
 346 static image and detecting interesting points using common feature detectors  
 347 such as SIFT (Lowe, 1999) or HOG (Dalal and Triggs, 2005). These features  
 348 are domain-specific for static images and are not easily extended to other  
 349 domains such as video (Le et al., 2011).

350 The approach taken by Stavens and Thrun (2010) learns its own domain-  
 351 optimized features instead of using pre-defined features, but still from static  
 352 images. A better approach to modeling videos is to learn image transitions  
 353 instead of working with static images. A Gated Restricted Boltzmann Ma-  
 354 chine (GRBM) (Memisevic and Hinton, 2007) has been used for this purpose  
 355 where the input,  $x$ , of the GRBM is the full image in one time frame and  
 356 the output  $y$  is the full image in the subsequent time frame. However, since  
 357 the network is fully connected to the image the method does not scale well

---

<sup>1</sup><http://www.nada.kth.se/cvap/actions/>

358 to larger images and local transformations at multiple locations must be  
359 re-learned.

360 A convolutional version of the GRBM using probabilistic max-pooling is  
361 presented by Taylor et al. (2010). The use of convolution reduces the number  
362 of parameters to learn, allows for larger input sizes, and better handles the  
363 local affine transformations that can appear anywhere in the image. The  
364 model was validated on synthetic data and a number of benchmark data  
365 sets, including the KTH activity recognition data set.

366 The work by Le et al. (2011) presents an unsupervised spatio-temporal  
367 feature learning method using an extension of Independent Subspace Analysis  
368 (ISA) (Hyvärinen et al., 2009). The extensions include hierarchical (stacked)  
369 convolutional ISA modules together with pooling. A disadvantage of ISA is  
370 that it does not scale well to large input sizes. The inclusion of convolution  
371 and stacking solves this problem by learning on smaller patches of input  
372 data. The method is validated on a number of benchmark sets, including  
373 KTH. One advantage of the method is that the use of ISA reduces the need  
374 for tweaking many of the hyperparameters seen in RBM-based methods, such  
375 as learning rate, weight decay, convergence parameters, etc.

376 Modeling temporal relations in video have also been done using temporal  
377 pooling. The work by Bo Chen and de Freitas (2010) uses convolutional  
378 RBMs as building blocks for spatial pooling and then performs temporal  
379 pooling on the spatial pooling units. The method is called Space-Time Deep  
380 Belief Network (ST-DBN). The ST-DBN allows for invariance and statistical  
381 dependencies in both space and time. The method achieved superior perfor-  
382 mance on applications such as action recognition and video denoising when

383 compared to a standard convolutional DBN.

384 The use of temporal coherence for modeling videos is done by Zou et al.  
385 (2011), where an auto-encoder with a L1-cost on the temporal difference on  
386 the pooling units is used to learn features that improve object recognition  
387 on still images. The work by Hyvärinen et al. (2003) also uses temporal  
388 information as a criterion for learning representations.

389 The use of deep learning, feature learning, and convolution with pooling  
390 has propelled the advances in video processing. Modeling streams of video is  
391 a natural continuation for deep learning algorithms since they have already  
392 been shown to be successful at building useful features from static images. By  
393 focusing on learning temporal features in videos, the performance on static  
394 images can be improved, which motivates the need for continuing developing  
395 deep learning algorithms that capture temporal relations. The early attempts  
396 at extending deep learning algorithms to video data was done by modeling the  
397 transition between two frames. The use of temporal pooling extends the time-  
398 dependencies a model can learn beyond a single frame transition. However,  
399 the time-dependency that has been modeled is still just a few frames. A  
400 possible future direction for video processing is to look at models that can  
401 learn longer time-dependencies.

#### 402 *4.2. Stock market prediction*

403 Stock market data are highly complex and difficult to predict, even for  
404 human experts, due to a number of external factors, e.g., politics, global  
405 economy, and trader expectation. The trends in stock market data tend to  
406 be nonlinear, uncertain, and non-stationary. Figure 7 shows the Dow Jones  
407 Industrial Average (DJOI) over a decade. According to the Efficient Market



Figure 7: Dow Jones Industrial Average (DJOI) over a period of 10 years.

408 Hypothesis (EMH) (Fama, 1965), stock market prices follow a random walk  
 409 pattern, meaning that a stock has the same probability to go up as it has  
 410 to go down, resulting in that predictions can not have more than 50% accu-  
 411 racy (Tsai and Hsiao, 2010). The EMH state that stock prices are largely  
 412 driven by "news" rather than present and past prices. However, it has also  
 413 been argued that stock market prices do not follow a random walk and that  
 414 they can be predicted (Malkiel, 2003). The landscape for acquiring both  
 415 news and stock information looks very different today than it did decades  
 416 ago. As an example, it has been shown that predicted stock prices can be  
 417 improved if further information is extracted from online social media, such  
 418 as Twitter feeds (Bollen et al., 2011) and online chat activity (Gruhl et al.,  
 419 2005).

420 One model that has emerged and shown to be suitable for stock market



421 prediction is the artificial neural network (ANN) (Atsalakis and Valavanis,  
422 2009). This is due to its ability to handle non-linear complex systems. A  
423 survey of ANNs applied to stock market prediction is given in Li and Ma  
424 (2010). However, most approaches of ANN applied to stock prediction have  
425 given unsatisfactory results (Agrawal et al., 2013). Neural networks with  
426 feedback have also been tried, such as recurrent versions of TDNN (Kim,  
427 1998), wavelet transformed features with an RNN (Hsieh et al., 2011), and  
428 echo state networks (Lin et al., 2009). Many of these methods are applied di-  
429 rectly on the raw data, while other papers focus more on the feature selection  
430 step (Tsai and Hsiao, 2010).

431 In summary, it can be concluded that there is still room to improve ex-  
432 isting techniques for making safe and accurate stock prediction systems. If  
433 additional information from sources that affect the stock market can be mea-  
434 sured and obtained, such as general public opinions from social media (Bollen  
435 et al., 2011), trading volume (Zhu et al., 2008), market specific domain knowl-  
436 edge, and political and economical factors, it can be combined together with  
437 the stock price data to achieve higher stock price predictions (Agrawal et al.,  
438 2013). The limited success of applying small, one layer neural networks for  
439 stock market prediction and the realization that there is a need to add more  
440 information to make better predictions indicate that a future direction for  
441 stock market prediction is to apply the combined data to more powerful  
442 models that are able to handle such complex, high-dimensional data. Deep  
443 learning methods for multivariate time-series fit this description and provide  
444 new interesting approach for the financial field and a new challenging appli-  
445 cation for the deep learning community, which to the authors knowledge has

446 not yet been tried.

### 447 4.3. Speech recognition

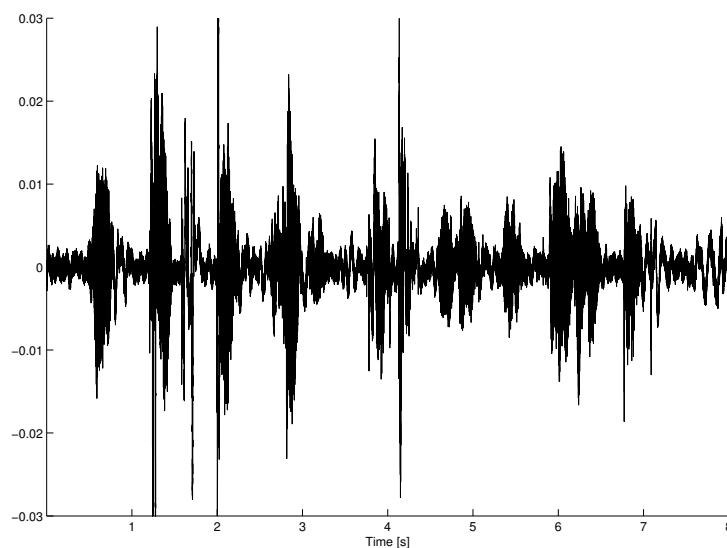


Figure 8: Raw acoustic signal of the utterance of the sentence "The quick brown fox jumps over the lazy dog".

448 Speech recognition is one area where deep learning has made significant  
449 progress (Hinton et al., 2012). The problem of speech recognition can be  
450 divided into a variety of sub-problems, such as speaker identification (Lee  
451 et al., 2009a), gender identification (Lee et al., 2009b; Parris and Carey,  
452 1996), speech-to-text (Furui et al., 2004) and acoustic modeling. The raw  
453 input data is single channel and highly time and frequency dependent, see  
454 Figure 8. A common approach is to use pre-set features that are designed  
455 for speech processing such as Mel-frequency cepstral coefficients (MFCC).

456 For decades, Hidden Markov Models (HMMs) (Rabiner and Juang, 1986)  
457 have been the state-of-the-art technique for speech recognition. A common

method for discretization of the input data for speech that is required by the HMM is to use Gaussian mixture models (GMM). More recently however, the Restricted Boltzmann Machines (RBM) have shown to be an adequate alternative for replacing the GMM in the discretization step. A classification error of 20.7% on the TIMIT speech recognition data set<sup>2</sup> was achieved by (Mohamed et al., 2012) by training a RBM on MFCC features. A similar setup has been used for large vocabulary speech recognition by Dahl et al. (2012). A convolutional deep belief networks was applied by Lee et al. (2009b) to audio data and evaluated on various audio classification tasks.

A number of variations on the RBM have also been tried on speech data. The mean-covariance RBM (mcRBM) (Ranzato and Hinton, 2010; Ranzato et al., 2010) achieved a classification error of 20.5% on the TIMIT data set by Dahl et al. (2010). A conditional RBM (cRBM) was modified by Mohamed and Hinton (2010) by including connections from future instead of only having connections from the past, which presumably gave better classification because the near future is more relevant than the more distant past.

Earlier, a Time-Delay Neural Network (TDNN) has been used for speech recognition (Waibel et al., 1989) and a review of TDNN architectures for speech recognition is given by Sugiyama et al. (1991). However, it has been suggested that convolution over the frequency instead of the time is better since the HMM on top models the temporal information.

The recent work by Graves et al. (2013) uses a deep Long Short-term Memory Recurrent Neural Network (RNN) (Hochreiter and Schmidhuber,

---

<sup>2</sup><http://www ldc.upenn.edu/Catalog/>

1997) to achieve a classification error of 17.7% on the TIMIT data set, which is the best result to date. One difference between the approaches of RBM-HMM and RNN is that the RNN can be used as an 'end-to-end' model because it replaces a combination of different techniques that are currently used in sequence modeling, such as the HMM. However, both these approaches still rely on pre-defined features as input.

While using features such as MFCCs that collapse high dimensional speech sound waves into low dimensional encodings have been successful in speech recognition systems, such low dimensional encodings may lose some relevant information. On the other hand, there are approaches that build their own features instead of using pre-defined features. The work by Jaitly and Hinton (2011) used raw speech as input to a RBM and achieved a classification error of 21.8% on the TIMIT data set. Another approach that uses raw data is learning the auditory codes using spiking population code (Smith and Lewicki, 2005). In this model, each spike encodes the precise time position and magnitude of a localized, time varying kernel function. The learned representations (basis vectors) show a striking resemblance to the cochlear filters in the auditory cortex.

Similarly sparse coding for audio classification is used by Grosse et al. (2007). The authors used features as input and a shift-invariant sparse coding model that reconstructs a time-series input using all the basis functions in all possible shifts. The model was evaluated on speaker identification and music genre classification.

A multimodal framework was explored by Ngiam et al. (2011) where video data of spoken digits and letters were combined with the audio data

507 to improve the classification.

508 In conclusion, there have been a lot of recent improvements to the pre-  
509 vious dominance of the features-GMM-HMM structure that has been used  
510 in speech recognition. First, there is a trend towards replacing GMM with  
511 a feature learning model such as deep belief networks or sparse coding. Sec-  
512 ond, there is a trend towards replacing HMM with other alternatives. One of  
513 them is the conditional random field (CRF) (Lafferty et al., 2001) that have  
514 been shown to outperform HMM, see for example the work by van Kasteren  
515 et al. (2008) and Bengio and Frasconi (1996). However, to date, the best  
516 reported result is replacing both parts of GMM-HMM with RNN (Graves  
517 et al., 2013). A next possible step for speech processing would be to replace  
518 the pre-made features with algorithms that build even better features from  
519 raw data.

#### 520 4.4. *Music recognition*

521 Music recognition is similar to speech recognition with the exception that  
522 the data can be multivariate and either presented as raw acoustic signals  
523 or by discrete chords. In music recognition, a number of sub-problems are  
524 considered, such as music annotation (genre, chord, instrument, mood classi-  
525 fication), music retrieval (text-based content search, content-based similarity  
526 retrieval, organization), and tempo identification. For music recognition,  
527 a commonly used set of features are MFCCs, chroma, constant-Q spectro-  
528 grams (CQT) (Schoerhuber and Klapuri, 2010), local contrast normalization  
529 (LCN) (LeCun et al., 2010), or Compressive Sampling (CS) (Chang et al.,  
530 2010). However, there is an increasing interest in learning the features from  
531 the data instead of using highly engineered features based on acoustic knowl-

edge. A widely used data set for music genre recognition is GTZAN<sup>3</sup>. Even though it is possible to solve many tasks on text-based meta-data, such as user data (playlists, song history, social structure), there is still a need for content-based analysis. The reasons for this is that manual labeling is inefficient due to the large amount of music content and some tasks require the well-trained ear of an expert, e.g., chord recognition.

The work by Humphrey et al. (2013) gives a review and future directions for music recognition. In this work, three deficiencies are identified: hand-crafted features are sub-optimal and unsustainable to develop for each task, shallow architectures are fundamentally limited, and short-time analysis cannot encode a musically meaningful structure. To handle these deficiencies it is proposed to learn features automatically, apply deep architectures, and model longer time-dependencies than the current use of data in milliseconds.

The work by Nam et al. (2012) addresses the first deficiency by presenting a processing pipeline for automatically learning features for music recognition. The model follows the structure of a high-dimensional single layer network with max-pooling separately after learning the features (Coates et al., 2010). The input data is taken from multiple audio frames and fed into three different feature learning algorithms, namely K-means clustering, sparse coding, and RBM. The learned features gave better performance compared to MFCC, regardless of the feature learning algorithm.

Sparse coding have been used by Grosse et al. (2007) for learning features for music genre recognition. The work by Henaff et al. (2011) used Predictive Sparse Decomposition (PSD), which is similar to sparse coding, and achieved

---

<sup>3</sup>[http://marsyas.info/download/data\\_sets](http://marsyas.info/download/data_sets)

556 an accuracy of 83.4% on the GTZAN data. In this work, the features are au-  
557 tomatically learned from CTQ spectrograms in an unsupervised manner. The  
558 learned features capture information about which chords are being played in  
559 a particular frame and produce comparable results to hand-crafted features  
560 for the task of genre recognition. A limitation, however, is that it ignores  
561 temporal dependencies between frames.

562 Convolutional DBNs were used by Lee et al. (2009b) to learn features from  
563 speech and music spectrograms and from engineered features by Dieleman  
564 et al. (2011). The work by (Hamel and Eck, 2010) also uses convolutional  
565 DBN to achieve an accuracy of 84.3% on the GTZAN dataset.

566 Self-taught learning have also been used for music genre classification.  
567 The self-taught learning framework attempts to use unlabeled data that does  
568 not share the labels of the classification task to improve classification perfor-  
569 mance (Raina et al., 2007; Jialin Pan and Yang, 2010). Self-taught learning  
570 and sparse coding are used by Markov and Matsui (2012) where unlabeled  
571 data from other music genres other than in the classification task was used  
572 to train the model.

573 In conclusion, there are many works that use unsupervised feature learn-  
574 ing methods for music recognition. The motivation for using deep networks  
575 is that music itself is structured hierarchically by a combination of chords,  
576 melodies and rhythms that creates motives, phrases, sections and finally en-  
577 tire pieces (Humphrey et al., 2013). Just like in speech recognition, the input  
578 data is often in some form of spectrograms. Many works leave the natural  
579 step of learning features from raw data as future work (Nam, 2012). Still, as  
580 proposed by (Humphrey et al., 2013), even though convolutional networks

581 have given good results on time-frequency representations of audio, there is  
 582 room for discovering new and better models.

#### 583 4.5. Motion capture data

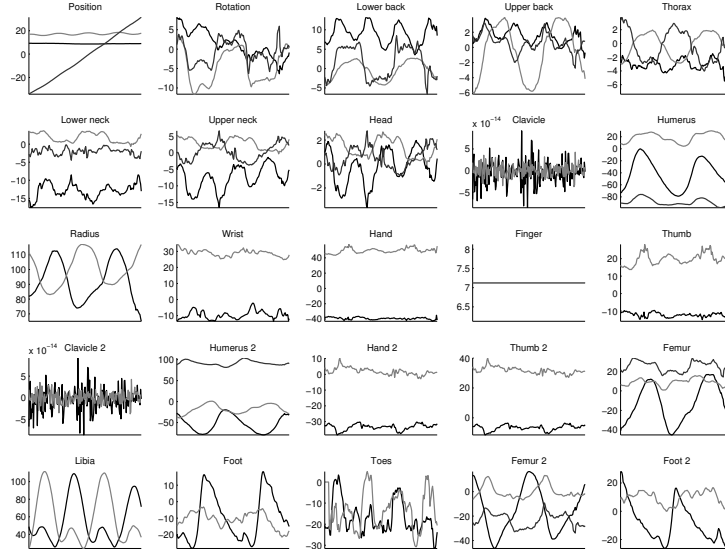


Figure 9: A sequence of human motion from the CMU motion capture data set.

584 Modeling human motion has several applications such as tracking, activ-  
 585 ity recognition, style and content separation, person identification, computer  
 586 animation, and synthesis of new motion data. Motion capture data is col-  
 587 lected from recordings of movements from several points on the body of a  
 588 human actor. These points can be captured by cameras that either track  
 589 the position of strategically placed markers (usually at joint centers) or uses  
 590 vision-based algorithms for tracking points of interest (Gleicher, 2000). The  
 591 points are represented as 3D Cartesian coordinates over time and are used to  
 592 form a skeletal structure with constant limb lengths by translating the points  
 593 to relative joint angles. The joint angles can be expressed in Euler angles,



594 4D quaternions, or exponential map parameterization (Grassia, 1998) and  
595 can have 1-3 degrees of freedom (DOF) each. The full data set consists of  
596 the orientation and translation of the root and all relative joint angles for  
597 each time frame as well as the constant skeleton model. The data is noisy,  
598 high-dimensional, and multivariate with complex nonlinear relationships. It  
599 has a lower frequency compared to speech and music data and some of the  
600 signals may be task-redundant.

601 Some of the traditional approaches include the work by Brand and Hertz-  
602 mann (2000), which models both the style and content of human motion using  
603 Hidden Markov Models (HMMs). The different styles were learned from un-  
604 labeled data and the trained model was used to synthesize motion data. A  
605 linear dynamical systems was used by Chiappa et al. (2009) to model three  
606 different motions of a human performing the task of holding a cup that has  
607 a ball attached to it with a string and then try to catch the ball into the cup  
608 (game of Balero). A Bayesian mixture of linear Gaussian state-space models  
609 (LGSSM) was trained with data from a human learner and used to generate  
610 new motions that was clustered and simulated on a robotic manipulator.

611 Both HMMs and linear dynamical systems are limited by their ability  
612 to model complex full-body motions. The work by Wang et al. (2007) uses  
613 Gaussian Processes to model three styles of locomotive motion (walk, run,  
614 stride) from the CMU motion capture data set<sup>4</sup>, see Figure 9. The CMU  
615 data set have also been used to generate motion capture from just a few  
616 initialization frames with a Temporal RBM (TRBM) (Sutskever and Hin-  
617 ton, 2006) and a conditional RBM (cRBM) Taylor et al. (2007). Better

---

<sup>4</sup><http://mocap.cs.cmu.edu/>

618 modeling and smoother transition between different styles of motions was  
619 achieved by adding a second hidden layer to the cRBM, using the Recurrent  
620 TRBM (Sutskever et al., 2008), and using the factored conditional RBM  
621 (fcRBM) (Taylor and Hinton, 2009). The work by Längkvist and Loutfi  
622 (2012) restructures an auto-encoder to resemble a cRBM but is used to per-  
623 form classification on the CMU motion capture data instead of generating  
624 new sequences. The drawbacks with general-purpose models such as Gaus-  
625 sian Processes and cRBM are that prior information about motion is not  
626 utilized and they have a costly approximation sampling procedure.

627 An unsupervised hierarchical model that is specifically designed for mod-  
628 eling locomotion styles was developed by Pan and Torresani (2009) and builds  
629 on the Hierarchical Bayesian Continuous Profile Model (HB-CPM). A Dy-  
630 namic Factor Graph (DFG), which is an extension of factor graphs, was  
631 introduced by Mirowski and LeCun (2009) and used on motion capture data  
632 to fill in missing data. The advantage of DFG is that it has a constant parti-  
633 tion function which avoids the costly approximation sampling procedure that  
634 is used in a cRBM.

635 In summary, analyzing and synthesizing motion capture data is a chal-  
636 lenging task and it encourages researchers to further improve learning algo-  
637 rithms for dealing with complex, multivariate time-series data. A motiva-  
638 tion for using deep learning algorithms for motion capture data is that it  
639 has been suggested that human motion is composed of elementary building  
640 blocks (motion templates) and any complex motion is constructed from a  
641 library of these previously learned motion templates (Flash and Hochner,  
642 2005). Deep networks can, in an unsupervised manner, learn these motion

643 templates from raw data and use them to form complex human motions.  
644 Motion capture data also provides an interesting platform for feature learn-  
645 ing from raw data since there is no commonly used feature set for motion  
646 capture data. Therefore, the success of applying deep learning algorithms to  
647 motion data can inspire learning features from raw data in other time-series  
648 problems as well.

#### 649 4.6. *Electronic nose data*

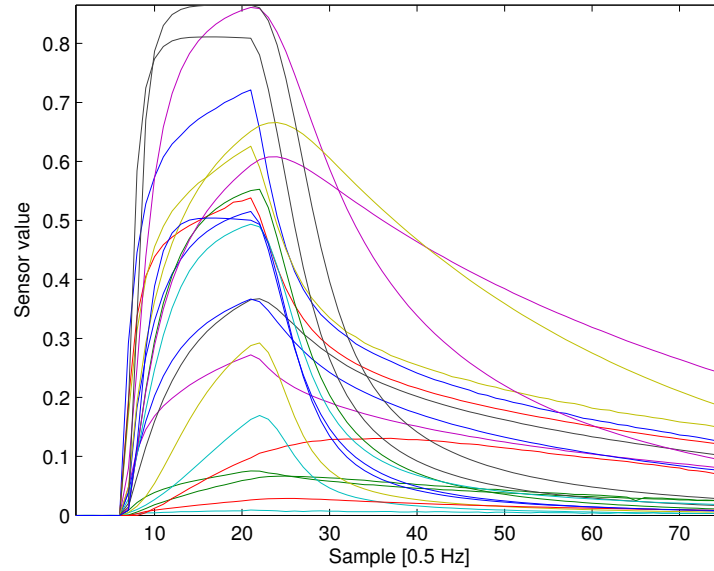


Figure 10: Normalized data from an array of electronic nose sensors.

650 Machine olfaction (Osuna et al., 2003; Gardner and Bartlett, 1999) is a  
651 field that seeks to quantify and analyze odours using an electronic nose (e-  
652 nose). An e-nose is composed of an array of selective gas sensors together  
653 with pattern recognition techniques. Figure 10 shows the readings from an e-  
654 nose sensor array. The number of sensors in the array typically ranges from  
655 4-30 sensors and are therefore, just like motion capture data, multivariate

656 and may contain redundant signals. The data is also unintuitive and there is  
657 a lack of expert knowledge that can guide the design of features. E-noses are  
658 mostly used in practice for industrial applications such as measuring food,  
659 beverage (Gardner et al., 2000b), and air quality (Zampolli et al., 2004), gas  
660 identification, and gas source localization Bennetts et al. (2011), but also has  
661 medical applications such as bacteria identification (Dutta et al., 2002) and  
662 diagnosis (Gardner et al., 2000a).

663 The traditional approach of analyzing e-nose data involves extracting  
664 information in the static and dynamic phases of the signals (Gutierrez-Osuna,  
665 2002) for the use of static pattern analysis techniques (PCA, discriminant  
666 function analysis, cluster analysis and neural networks). Some commonly  
667 used features are the static sensor response, transient derivatives (Trincavelli  
668 et al., 2010), area under the curve (Carmona et al., 2006), model parameter  
669 identification (Vembu et al., 2012), and dynamic analysis (Hines et al., 1999).

670 A popular approach for modeling e-nose data is the Time-Delay Neural  
671 Networks (TDNN) (Waibel et al., 1989). It has been used for identifying  
672 the smell of spices (Zhang et al., 2003), ternary mixtures (Vito et al., 2007),  
673 optimum fermentation time for black tea (Bhattacharya et al., 2008), and  
674 vintages of wine (Yamazaki et al., 2001). An RNN have been used for odour  
675 localization with a mobile robot (Duckett et al., 2001).

676 The work by Vembu et al. (2012) compares the gas discrimination and  
677 localization between three approaches: SVM on raw data, SVM on features  
678 extracted from auto-regressive and linear dynamical systems, and finally a  
679 SVMs with kernels specialized for structured data (Gärtner, 2003). The SVM  
680 with built-in time-aware kernels performed better than techniques that used

681 feature extraction, even though the features captured temporal information.

682 More recently, an auto-encoder, RBM, and cRBM have been used for  
683 bacteria identification (Långkvist and Loutfi, 2011) and fast classification of  
684 meat spoilage markers (Långkvist et al., 2013).

685 E-nose data introduces the challenge of improving models that can deal  
686 with redundant signals. It is not feasible to produce tailor-made sensors for  
687 each possible individual gas and combinations of gases of interest. Therefore  
688 the common approach is to use an array of sensors with different properties  
689 and leave the discrimination to the pattern analysis software. It is also not  
690 desirable to construct new feature sets for each e-nose application so a data-  
691 driven feature learning method is useful. The early works on e-nose data  
692 create feature vectors of simple features for each signal such as the static  
693 response or the slope of dynamic response and then feed it to a classifier.  
694 Recently, the use of dynamic models such as neural networks with tapped  
695 delays and SVMs with kernels for structured data have shown to improve the  
696 performance over static approaches. The next step is to continue this trend  
697 of using dynamical models that constructs robust features that can deal with  
698 noisy inputs in order to quantify and classify odors in more challenging open  
699 environments with many different simultaneous gas sources.

#### 700 4.7. *Physiological data*

701 With physiological data we consider recordings such as electroencephalog-  
702 raphy (EEG), magnetoencephalography (MEG), electrocardiography (ECG),  
703 and wearable sensors for health monitoring. Figure 11 shows an example of  
704 how physiological data look like. The data can exist both as singular or  
705 multiple channels. The use of a feature learning algorithm is particularly

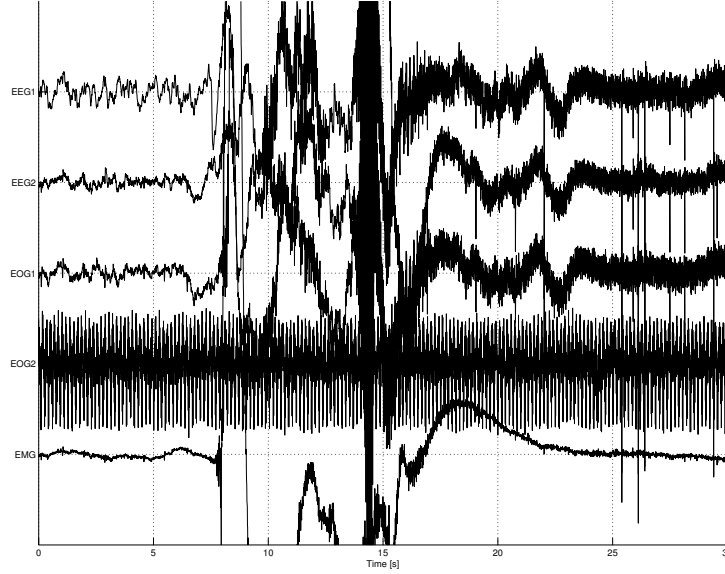


Figure 11: Data from EEG (top two signals), EOG (third and fourth signal), and EMG (bottom signal), recorded with a polysomnograph during sleep.

706 beneficial in medical applications because acquiring a labeled medical data  
 707 set is expensive since the data sets are often very large and require the la-  
 708 beling of an expert in the field.

709 The work by Mirowski et al. (2008) compares convolutional networks  
 710 with logistic regression and SVMs for epileptic seizure prediction from in-  
 711 tracranial EEG signals. The features that are used are hand-engineered bi-  
 712 variate features between channels that encode relationship between pairs of  
 713 EEG channels. The result was that convolutional networks achieved only 1  
 714 false-alarm prediction from 21 patients while the SVM had 10 false-alarms.  
 715 TDNN and ICA has also been used for EEG-based prediction of epileptic  
 716 seizures (Mirowski et al., 2007). The application of self-organizing maps  
 717 (SOM) to analyze EMG data is presented by Tucker (1999).

718 A RBM-based method that builds features from raw data for sleep stage

719 classification from 4-channel polysomnography data has been proposed by Långkvist  
720 et al. (2012). A similar setup was used by Wulsin et al. (2011) for model-  
721 ing single channel EEG waveforms used for anomaly detection. A DBN  
722 is used by (Wang and Shang, 2013) to automatically extract features from  
723 raw unlabelled physiological data and achieves better classification than a  
724 feature-based approach. These recent works show that DBNs can be applied  
725 to raw physiological data to effectively learn relevant features.

726 A source separation method tailor-made to EEG and MEG signals is pro-  
727 posed by Hyvärinen et al. (2010). The data is preprocessed by short-time  
728 Fourier transforms and then fed to an ICA. The work shows that tempo-  
729 ral correlations are adequately taken into account. Independent Component  
730 Analysis (ICA) has provided to be a new tool to analyze time series and is  
731 a unifying framework that combines sparseness, temporal coherence, topog-  
732 raphy and complex cell pooling in a single model (Hyvärinen et al., 2003).  
733 A method for how to order the independent components for time-series is  
734 explored by Cheung and Xu (2001).

735 Self-taught learning has been used with time-series data from wearable  
736 hand-motion sensors (Amft, 2011).

737 The field of physiological data is large and many different methods have  
738 been used. The characteristics of physiological data could be particularly  
739 interesting for the deep learning community because it can be used to explore  
740 the feasibility of learning features from raw data, which hopefully can inspire  
741 similar approaches in other time-series domains.

Table 2: A summary of commonly used time-series problems.

Problem	Multi-variate	Raw data	Frequency rich	Common features	Common method	Benchmark set
Stock prediction	-	✓	-	-	ANN	DJIA
Video	✓	✓	-	SIFT, HOG	ConvRBM	KTH
Speech Recognition	-	(✓)	✓	MFCC	RBM, RNN	TIMIT
Music recognition	✓	-	✓	Chroma, MFCC	ConvRBM	GTZAN
Motion capture	✓	✓	-	-	cRBM	CMU
E-nose	✓	✓	-	Many	TDNN	-
Physiological data	✓	(✓)	✓	Many, spectrogram	RBM, AE	PhysioNET

#### 742 4.8. Summary

743 Table 2 gives a summary of the time-series problems that have been pre-  
744 sented in this section. The first column indicates if the data is multivariate  
745 (or only contains one signal, univariate). Stock prediction is often viewed as  
746 a single channel problem, which explains the difficulties to produce accurate  
747 prediction systems, since stocks depend on a myriad of other factors, and  
748 arguably not at all on past values of the stock itself. For speech recognition,  
749 the use of multimodal sources can improve performance (Ngiam et al., 2011).

750 The second column shows which problems have attempted to create fea-  
751 tures purely from raw data. Only a few works have attempted this with  
752 speech recognition (Jaitly and Hinton, 2011; Smith and Lewicki, 2005) and  
753 physiological data (Wulsin et al., 2011; Långkvist et al., 2012; Wang and  
754 Shang, 2013). To the authors knowledge, learning features from raw data  
755 has not been attempted in music recognition. The process of constructing  
756 features from raw data has been well demonstrated for vision-tasks but is



cautiously used for time-series problems. Models such as TDNN, cRBM and convolutional RBMs are well suited for being applied to raw data (or slightly pre-processed data).

The third column indicates which time-series problems have valuable information in the frequency-domain. For frequency-rich problems, it is uncommon to attempt to learn features from raw data. A reason for this is that current feature learning algorithms are yet not well-suited for learning features in the frequency-domain.

The fourth column displays some common features that have been used in the literature. SIFT and HOG have been applied to videos even though those features are developed for static images. Chroma and MFCC have been applied to music recognition, even though they are developed for speech recognition. The e-nose community have tried a plethora of features. E-nose data is a relatively new field where a hand-crafted feature set have not been developed since this kind of data is complex and unintuitive. For physiological data, the used features are often a combination of application-specific features from previous works or hand-crafted features.

The fifth column reports the most commonly used method(s), or current state-of-the-art, for each time-series problem. For stock prediction, the progress has stopped at classical neural networks. The current state-of-the-art augments additional information beside the stock data. For high-dimensional temporal data such as video and music recognition, the convolutional version of RBM have been successful. In recent years, the RBM have been used for speech recognition but the current state-of-the-art is achieved with an RNN. The cRBM introduced motion capture data to the deep learn-

ing community and it is an interesting problem to explore with other methods. Single layer neural networks with temporal capabilities have been used to model e-nose data and the use of deep networks is an interesting future direction for modeling e-nose data.

And finally, the last column indicates a typical benchmark set for each problem. There is currently no well-known publicly available benchmark data set for e-nose data. For deep learning to enter the field of e-nose data it requires a large, well-organized data set that would benefit both communities. A data base of physiological data is available from PhysioNET (Goldberger et al., 2000 (June 13)).

## 5. Conclusion

Unsupervised feature learning and deep learning techniques have been successfully applied to a variety of domains. While much focus in deep learning and unsupervised feature learning have been in the computer vision domain, this paper has reviewed some of the successful applications of deep learning methods to the time-series domain. Some of these approaches have treated the input as static data but the most successful ones are those that have modified the deep learning models to better handle time-series data.

The problem with processing time-series data as static input is that the importance of time is not captured. Modeling time-series faces many of the same challenges as modeling static data, such as coping with high-dimensional observations and nonlinear relationships between variables, however, by simply ignoring time and applying models of static data to time series

806 one disregards much of the rich structure present in the data. When taking  
807 this approach, the context of the current input frame is lost and the only  
808 time-dependencies that are captured is within the input size. In order to  
809 capture long-term dependencies, the input size has to be increased, which  
810 can be impractical for multivariate signals or if the data has very long-term  
811 dependencies. The solution is to use a model that incorporates temporal  
812 coherence, performs temporal pooling, or models sequences of hidden unit  
813 activations.

814     The choice of model and how the data should be presented to the model  
815 is highly dependent on the type of data. Within a chosen model there are  
816 additional design choices in terms of connectivity, architecture, and hyperpa-  
817 rameters. For these reasons, even though many unsupervised feature learning  
818 models offer to relieve the user of having to come up with useful features for  
819 the current domain, there are still many challenges for applying them to time-  
820 series data. It is also worth noting that many works that construct useful  
821 features from the input data actually still use input data from pre-processed  
822 features.

823     Deep learning methods offer better representation and classification on a  
824 multitude of time-series problems compared to shallow approaches when con-  
825 figured and trained properly. There is still room for improving the learning  
826 algorithms specifically for time-series data, e.g., performing signal selection  
827 that deals with redundant signals in multivariate input data. Another possi-  
828 ble future direction is to develop models that change their internal architec-  
829 ture during learning or use model averaging in order to capture both short  
830 and long-term time dependencies. Further research in this area is needed to

831 develop algorithms for time-series modeling that learn even better features  
832 and are easier and faster to train. Therefore, there is a need to focus less on  
833 the pre-processing pipeline for a specific time-series problem and focus more  
834 on learning better feature representations for a general-purpose algorithm for  
835 structured data, regardless of the application.

## 836 References

- 837 Agrawal, J.G., Chourasia, V.S., Mittra, A.K., 2013. State-of-the-art in stock  
838 prediction techniques. *International Journal of Advanced Research in Elec-*  
839 *trical, Electronics and Instrumentation Engineering* 2, 1360–1366.
- 840 Amft, O., 2011. Self-taught learning for activity spotting in on-body mo-  
841 tion sensor data, in: *ISWC 2011: Proceedings of the IEEE International*  
842 *Symposium on Wearable Computing*, IEEE. pp. 83–86.
- 843 Atsalakis, G.S., Valavanis, K.P., 2009. Surveying stock market forecasting  
844 techniques Íc part ii: Soft computing methods. *Expert Systems with Ap-*  
845 *plications* 36, 5932 – 5941.
- 846 Bengio, Y., 2007. Learning deep architectures for AI. Technical Report 1312.  
847 Dept. IRO, Universite de Montreal.
- 848 Bengio, Y., Courville, A., Vincent, P., 2012. Unsupervised Feature Learning  
849 and Deep Learning: A Review and New Perspectives. Technical Report  
850 arXiv:1206.5538. U. Montreal. URL: <http://arxiv.org/abs/1206.5538>.
- 851 Bengio, Y., Frasconi, P., 1996. Input-output HMM’s for sequence processing.  
852 *IEEE Transactions on Neural Networks* 7(5), 1231–1249.

- 853 Bengio, Y., Lamblin, P., Popovici, D., Larochelle, H., 2007. Greedy layer-  
854 wise training of deep networks. *Advances in neural information processing*  
855 systems 19, 153.
- 856 Bengio, Y., LeCun, Y., 2007. Scaling learning algorithms towards AI, in:  
857 Bottou, L., Chapelle, O., DeCoste, D., Weston, J. (Eds.), *Large-Scale*  
858 *Kernel Machines*, MIT Press.
- 859 Bengio, Y., Simard, P., Frasconi, P., 1994. Learning longterm dependencies  
860 with gradient descent is difficult. *IEEE Transactions on Neural Networks*  
861 5(2), 157–166.
- 862 Bengio, Y., Yao, L., Alain, G., Vincent, P., 2013. Generalized denoising  
863 auto-encoders as generative models. *CoRR* abs/1305.6663.
- 864 Bennetts, V.H., Lilienthal, A.J., Neumann, P.P., Trincavelli, M., 2011. Mo-  
865 bile robots for localizing gas emission sources on landfill sites: is bio-  
866 inspiration the way to go? *Frontiers in neuroengineering* 4.
- 867 Bhattacharya, N., Tudu, B., Jana, A., Ghosh, D., Bandhopadhyaya, R.,  
868 Bhuyan, M., 2008. Preemptive identification of optimum fermentation time  
869 for black tea using electronic nose. *Sensors and Actuators B: Chemical* 131,  
870 110–116.
- 871 Bo Chen, Jo-Anne Ting, B.M., de Freitas, N., 2010. Deep learning of in-  
872 variant spatio-temporal features from video, in: *NIPS 2010 Deep Learning*  
873 *and Unsupervised Feature Learning Workshop*.
- 874 Bollen, J., Mao, H., Zeng, X., 2011. Twitter mood predicts the stock market.  
875 *Journal of Computational Science* 2, 1 – 8.

- 876 Bottou, L., 2010. Large-scale machine learning with stochastic gradient de-  
 877 scent, in: Lechevallier, Y., Saporta, G. (Eds.), Proceedings of the 19th In-  
 878 ternational Conference on Computational Statistics (COMPSTAT'2010),  
 879 Springer, Paris, France. pp. 177–187. URL: [http://leon.bottou.org/](http://leon.bottou.org/papers/bottou-2010)  
 880 [papers/bottou-2010](http://leon.bottou.org/papers/bottou-2010).
- 881 Brand, M., Hertzmann, A., 2000. Style machines, in: Proceedings of the 27th  
 882 annual conference on Computer graphics and interactive techniques, ACM  
 883 Press/Addison-Wesley Publishing Co., New York, NY, USA. pp. 183–192.
- 884 Carmona, M., Martinez, J., Zalacain, A., Rodriguez-Mendez, M.L., de Saja,  
 885 J.A., Alonso, G.L., 2006. Analysis of saffron volatile fraction by td-gc-ms  
 886 and e-nose. European Food Research and Technology 223, 96–101.
- 887 Chang, K., Jang, J., Iliopoulos, C., 2010. Music genre classification via com-  
 888 pressive sampling, in: Proceedings of the 11th International Conference on  
 889 Music Information Retrieval (ISMIR), pp. 387–392.
- 890 Cheung, Y., Xu, L., 2001. Independent component ordering in ica time series  
 891 analysis. Neurocomputing 41, 145–152.
- 892 Chiappa, S., Kober, J., Peters, J., 2009. Using bayesian dynamical systems  
 893 for motion template libraries. In Adv. in Neural Inform. Proc. Systems 21,  
 894 297–304.
- 895 Coates, A., Lee, H., Ng, A.Y., 2010. An Analysis of Single-Layer Networks  
 896 in Unsupervised Feature Learning. Engineering , 1–9.
- 897 Dahl, G., Yu, D., Deng, L., Acero, A., 2012. Context-dependent pre-  
 898 trained deep neural networks for large-vocabulary speech recognition. Au-

899     dio, Speech, and Language Processing, IEEE Transactions on 20, 30–42.  
900     doi:10.1109/TASL.2011.2134090.

901     Dahl, G.E., Ranzato, M., Mohamed, A., Hinton, G., 2010. Phone recogni-  
902     tion with the mean-covariance restricted boltzmann machine. Advances in  
903     Neural Information Processing Systems 23, 469–477.

904     Dalal, N., Triggs, B., 2005. Histograms of oriented gradients for human  
905     detection, in: In CVPR.

906     Dieleman, S., Brakel, P., Schrauwen, B., 2011. Audio-based music classifi-  
907     cation with a pretrained convolutional network, in: In The International  
908     Society for Music Information Retrieval (ISMIR).

909     Dietterich, T.G., 2002. Machine learning for sequential data: A review,  
910     in: Structural, Syntactic, and Statistical Pattern Recognition, Springer-  
911     Verlag. pp. 15–30.

912     Duckett, T., Axelsson, M., Saffiotti, A., 2001. Learning to locate an odour  
913     source with a mobile robot, in: Robotics and Automation, 2001. Proceed-  
914     ings 2001 ICRA. IEEE International Conference on, pp. 4017–4022 vol.4.  
915     doi:10.1109/ROBOT.2001.933245.

916     Dutta, R., Hines, E., Gardner, J., Boilot, P., 2002. Bacteria classification  
917     using cyranose 320 electronic nose. Biomedical Engineering Online 1, 4.

918     Erhan, D., Bengio, Y., Courville, A., Manzagol, P., Vincent, P., Bengio, S.,  
919     2010. Why does unsupervised pre-training help deep learning? Journal of  
920     Machine Learning Research 11, 625–660.

- 921 Fama, E.F., 1965. The behavior of stock-market prices. The Journal of  
922 Business 1, 34–105.
- 923 Flash, T., Hochner, B., 2005. Motor primitives in vertebrates and inverte-  
924 brates. Current Opinion in Neurobiology 15(6), 660–666.
- 925 Furui, S., Kikuchi, T., Shinnaka, Y., Hori, C., 2004. Speech-to-text and  
926 speech-to-speech summarization of spontaneous speech. Speech and Audio  
927 Processing, IEEE Transactions on 12, 401–408.
- 928 Gardner, J., Bartlett, P., 1999. Electronic Noses, Principles and Applications.  
929 Oxford University Press, New York, NY, USA.
- 930 Gardner, J.W., Shin, H.W., Hines, E.L., 2000a. An electronic nose system  
931 to diagnose illness. Sensors and Actuators B: Chemical 70, 19–24.
- 932 Gardner, J.W., Shin, H.W., Hines, E.L., Dow, C.S., 2000b. An electronic nose  
933 system for monitoring the quality of potable water. Sensors and Actuators  
934 B: Chemical 69, 336–341.
- 935 Gärtner, T., 2003. A survey of kernels for structured data. SIGKDD Explor.  
936 Newsl. 5, 49–58.
- 937 Gers, F.A., Schmidhuber, J., Cummins, F., 2000. Learning to Forget: Con-  
938 tinual Prediction with LSTM. Neural Computation 12, 2451–2471.
- 939 Gleicher, M., 2000. Animation from observation: Motion capture and motion  
940 editing. SIGGRAPH Computer Graphics 33, 51–54.



941 Goldberger, A.L., Amaral, L.A.N., Glass, L., Hausdorff, J.M., Ivanov,  
942 P.C., Mark, R.G., Mietus, J.E., Moody, G.B., Peng, C.K., Stan-  
943 ley, H.E., 2000 (June 13). PhysioBank, PhysioToolkit, and Phys-  
944 ioNet: Components of a new research resource for complex physio-  
945 logic signals. *Circulation* 101, e215–e220. *Circulation Electronic Pages*:  
946 <http://circ.ahajournals.org/cgi/content/full/101/23/e215>.

947 Grassia, F.S., 1998. Practical parameterization of rotations using the expo-  
948 nential map. *J. Graph. Tools* 3, 29–48.

949 Graves, A., Mohamed, A., Hinton, G., 2013. Speech recognition with deep re-  
950 current neural networks, in: *The 38th International Conference on Acous-*  
951 *tics, Speech, and Signal Processing (ICASSP)*.

952 Grosse, R., Raina, R., Kwong, H., Ng, A.Y., 2007. Shift-invariant sparse  
953 coding for audio classification, in: *Conference on Uncertainty in Artificial*  
954 *Intelligence (UAI)*.

955 Gruhl, D., Guha, R., Kumar, R., Novak, J., Tomkins, A., 2005. The predic-  
956 tive power of online chatter, in: *Proceedings of the eleventh ACM SIGKDD*  
957 *international conference on Knowledge discovery in data mining*, pp. 78–  
958 87.

959 Gutierrez-Osuna, R., 2002. Pattern analysis for machine olfaction: A review.  
960 *IEEE Sensors Journal* 2(3), 189–202.

961 Hamel, P., Eck, D., 2010. Learning features from music audio with deep belief  
962 networks, in: *11th International Society for Music Information Retrieval*  
963 *Conference (ISMIR)*.

- 964 Henaff, M., Jarrett, K., Kavukcuoglu, K., LeCun, Y., 2011. Unsupervised  
965 learning of sparse features for scalable audio classification, in: Proceedings  
966 of International Symposium on Music Information Retrieval (ISMIR'11).
- 967 Hines, E., Llobet, E., Gardner, J., 1999. Electronic noses: a review of signal  
968 processing techniques. *Circuits, Devices and Systems*, IEE Proceedings -  
969 146, 297–310.
- 970 Hinton, G., Deng, L., Yu, D., Dahl, G.E., Mohamed, A.r., Jaitly, N., Senior,  
971 A., Vanhoucke, V., Nguyen, P., Sainath, T.N., et al., 2012. Deep neural  
972 networks for acoustic modeling in speech recognition: The shared views of  
973 four research groups. *Signal Processing Magazine*, IEEE 29, 82–97.
- 974 Hinton, G., Salakhutdinov, R., 2006. Reducing the dimensionality of data  
975 with neural networks. *Science* 313(5786), 504–507.
- 976 Hinton, G.E., 2002. Training products of experts by minimizing contrastive  
977 divergence. *Neural Computation* 14, 1771 – 1800.
- 978 Hinton, G.E., 2012. A practical guide to training restricted boltzmann ma-  
979 chines, in: Montavon, G., Orr, G.B., Müller, K.R. (Eds.), *Neural Networks:*  
980 *Tricks of the Trade*. Springer Berlin Heidelberg. volume 7700 of *Lecture*  
981 *Notes in Computer Science*, pp. 599–619. URL: [http://dx.doi.org/10.](http://dx.doi.org/10.1007/978-3-642-35289-8_32)  
982 [1007/978-3-642-35289-8\\_32](http://dx.doi.org/10.1007/978-3-642-35289-8_32), doi:10.1007/978-3-642-35289-8\_32.
- 983 Hinton, G.E., Krizhevsky, A., Wang, S.D., 2011. Transforming auto-  
984 encoders, in: Proceedings of the 21th international conference on Artificial  
985 neural networks - Volume Part I, pp. 44–51.

986 Hinton, G.E., S., O., Y., T., 2006. A fast learning algorithm for deep belief  
987 nets. *Neural Computation* 18 , 1527–1554.

988 Hochreiter, S., Schmidhuber, J., 1997. Long short-term memory. *Neural*  
989 *Computation* 9, 1735–1780.

990 Hsieh, T.J., Hsiao, H.F., Yeh, W.C., 2011. Forecasting stock markets using  
991 wavelet transforms and recurrent neural networks: An integrated system  
992 based on artificial bee colony algorithm. *Applied Soft Computing* 11, 2510  
993 – 2525.

994 Humphrey, E.J., Bello, J.P., LeCun, Y., 2013. Feature learning and deep  
995 architectures: new directions for music informatics. *Journal of Intelligent*  
996 *Information Systems* 41, 461–481.

997 Hüsken, M., Stagge, P., 2003. Recurrent Neural Networks for Time Series  
998 Classification. *Neurocomputing* 50, 223–235.

999 Hyvärinen, A., Hurri, J., Vährynen, J., 2003. Bubbles: a unifying framework  
1000 for low-level statistical properties of natural image sequences. *J. Opt. Soc.*  
1001 *Am. A* 20, 1237–1252.

1002 Hyvärinen, A., Ramkumar, P., Parkkonen, L., Hari, R., 2010. Indepen-  
1003 dent component analysis of short-time Fourier transforms for spontaneous  
1004 EEG/MEG analysis. *NeuroImage* 49(1), 257–271.

1005 Hyvärinen, A., Hurri, J., Hoyer, P.O., 2009. *Natural Image Statistics*. vol-  
1006 ume 39. Springer.

1007 Jaitly, N., Hinton, G., 2011. Learning a better representation of speech  
1008 soundwaves using restricted boltzmann machines, in: Acoustics, Speech  
1009 and Signal Processing (ICASSP), 2011 IEEE International Conference on,  
1010 IEEE. pp. 5884–5887.

1011 Jialin Pan, S., Yang, Q., 2010. A survey on transfer learning. IEEE Trans-  
1012 actions On Knowledge and Data Engineering 22.

1013 Kamyshanska, H., Memisevic, R., 2013. On autoencoder scoring, in: Pro-  
1014 ceedings of the 30th International Conference on Machine Learning (ICML-  
1015 13), JMLR Workshop and Conference Proceedings. pp. 720–728.

1016 van Kasteren, T., Noulas, A., Kröse, B., 2008. Conditional random fields  
1017 versus hidden markov models for activity recognition in temporal sensor  
1018 data, in: In Proceedings of the 14th Annual Conference of the Advanced  
1019 School for Computing and Imaging (ASCI'08), The Netherlands.

1020 Kavukcuoglu, K., Ranzato, M., Fergus, R., Le-Cun, Y., 2009. Learning  
1021 invariant features through topographic filter maps, in: Computer Vision  
1022 and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on, IEEE.  
1023 pp. 1605–1612.

1024 Keogh, E., Kasetty, S., 2002. On the need for time series data mining bench-  
1025 marks: A survey and empirical demonstration, in: In proceedings of the  
1026 8th ACM SIGKDD International Conference on Knowledge Discovery and  
1027 Data Mining, pp. 102–111.

1028 Kim, S.S., 1998. Time-delay recurrent neural network for temporal correla-  
1029 tions and prediction. Neurocomputing 20, 253 – 263.

- 1030 Lafferty, J.D., McCallum, A., Pereira, F.C.N., 2001. Conditional random  
1031 fields: Probabilistic models for segmenting and labeling sequence data,  
1032 in: Proceedings of the Eighteenth International Conference on Machine  
1033 Learning, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA.  
1034 pp. 282–289.
- 1035 Längkvist, M., Coradeschi, S., Loutfi, A., Rayappan, J.B.B., 2013. Fast  
1036 Classification of Meat Spoilage Markers Using Nanostructured ZnO Thin  
1037 Films and Unsupervised Feature Learning. *Sensors* 13(2), 1578–1592.  
1038 Doi:10.3390/s130201578.
- 1039 Längkvist, M., Karlsson, L., Loutfi, A., 2012. Sleep stage classification using  
1040 unsupervised feature learning. *Advances in Artificial Neural Systems* 2012.  
1041 Doi:10.1155/2012/107046.
- 1042 Längkvist, M., Loutfi, A., 2011. Unsupervised feature learning for electronic  
1043 nose data applied to bacteria identification in blood, in: NIPS workshop  
1044 on Deep Learning and Unsupervised Feature Learning.
- 1045 Längkvist, M., Loutfi, A., 2012. Not all signals are created equal: Dynamic  
1046 objective auto-encoder for multivariate data, in: NIPS workshop on Deep  
1047 Learning and Unsupervised Feature Learning.
- 1048 Le, Q.V., Zou, W.Y., Yeung, S.Y., Ng, A.Y., 2011. Learning hierarchical  
1049 invariant spatio-temporal features for action recognition with independent  
1050 subspace analysis, in: Computer Vision and Pattern Recognition (CVPR).
- 1051 Le Roux, N., Bengio, Y., 2008. Representational power of restricted Boltz-

mann machines and deep belief networks. *Neural Computation* 20, 1631–1649.

LeCun, Y., Kavukvuoglu, K., Farabet, C., 2010. Convolutional networks and applications in vision, in: *Proc. International Symposium on Circuits and Systems (ISCAS’10)*, IEEE.

Lee, H., Ekanadham, C., Ng, A.Y., 2008. Sparse deep belief net model for visual area V2, in: *Advances in Neural Information Processing Systems* 20, pp. 873–880.

Lee, H., Grosse, R., Ranganath, R., Ng, A.Y., 2009a. Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations, in: *Twenty-Sixth International Conference on Machine Learning*.

Lee, H., Largman, Y., Pham, P., Ng, A.Y., 2009b. Unsupervised feature learning for audio classification using convolutional deep belief networks, in: *Advances in Neural Information Processing Systems* 22, pp. 1096–1104.

Li, Y., Ma, W., 2010. Applications of artificial neural networks in financial economics: A survey, in: *Proceedings of the 2010 International Symposium on Computational Intelligence and Design - Volume 01*, IEEE Computer Society. pp. 211–214.

Lin, X., Yang, Z., Song, Y., 2009. Short-term stock price prediction based on echo state networks. *Expert Systems with Applications* 36, 7313 – 7317.

Lowe, D., 1999. Object recognition from local scale-invariant features, in: *ICCV*.

- 1074 Luenberger, D., 1979. Introduction to Dynamic Systems: Theory, Models,  
1075 and Applications. Wiley.
- 1076 Lütkepohl, H., 2005. New Introduction to Multiple Time Series Analysis.  
1077 Springer-Verlag.
- 1078 Malkiel, B., 2003. The efficient market hypothesis and its critics. The Journal  
1079 of Economic Perspectives 17. [Http://dx.doi.org/10.2307/3216840](http://dx.doi.org/10.2307/3216840).
- 1080 Markov, K., Matsui, T., 2012. Music genre classification using self-taught  
1081 learning via sparse coding, in: Acoustics, Speech and Signal Processing  
1082 (ICASSP), 2012 IEEE International Conference on, pp. 1929–1932.
- 1083 Martens, J., Sutskever, I., 2012. Training deep and recurrent neural networks  
1084 with hessian-free optimization, in: Neural Networks: Tricks of the Trade.  
1085 Springer Berlin Heidelberg. volume 7700 of *Lecture Notes in Computer*  
1086 *Science*.
- 1087 Masci, J., Meier, U., Cireşan, D., Schmidhuber, J., 2011. Stacked convolu-  
1088 tional auto-encoders for hierarchical feature extraction, in: Proceedings of  
1089 the 21th international conference on Artificial neural networks - Volume  
1090 Part I, pp. 52–59.
- 1091 Memisevic, R., Hinton, G., 2007. Unsupervised learning of image transforma-  
1092 tions, in: IEEE Conference on Computer Vision and Pattern Recognition  
1093 (CVPR), pp. 1–8.
- 1094 Memisevic, R., Hinton, G.E., 2010. Learning to represent spatial transforma-  
1095 tions with factored higher-order boltzmann machines. Neural Computation  
1096 22, 1473–1492.

1097 Mirowski, P., LeCun, Y., 2009. Dynamic factor graphs for time series model-  
1098 ing. *Machine Learning and Knowledge Discovery in Databases* , 128–143.

1099 Mirowski, P., Madhavan, D., LeCun, Y., 2007. Time-delay neural networks  
1100 and independent component analysis for eeg-based prediction of epileptic  
1101 seizures propagation, in: *Association for the Advancement of Artificial*  
1102 *Intelligence Conference*.

1103 Mirowski, P.W., LeCun, Y., Madhavan, D., Kuzniecky, R., 2008. Comparing  
1104 SVM and convolutional networks for epileptic seizure prediction from in-  
1105 tracranial EEG, in: *Machine Learning for Signal Processing, 2008. MLSP*  
1106 *2008. IEEE Workshop on, IEEE*. pp. 244–249.

1107 Mohamed, A., Dahl, G.E., Hinton, G., 2012. Acoustic modeling using deep  
1108 belief networks. *IEEE Transactions on Audio, Speech, and Language Pro-*  
1109 *cessing archive* 20(1), 14–22.

1110 Mohamed, A., Hinton, G., 2010. Phone recognition using restricted boltz-  
1111 mann machines, in: *Acoustics Speech and Signal Processing (ICASSP),*  
1112 *2010 IEEE International Conference on*, pp. 4354–4357. doi:10.1109/  
1113 *ICASSP.2010.5495651*.

1114 Nam, J., 2012. Learning Feature Representations for Music Classification.  
1115 Ph.D. thesis. Stanford University.

1116 Nam, J., Herrera, J., Slaney, M., Smith, J.O., 2012. Learning Sparse Feature  
1117 Representations for Music Annotation and Retrieval, in: *In The Interna-*  
1118 *tional Society for Music Information Retrieval (ISMIR)*, pp. 565–570.



- 1119 Nanopoulos, A., Alcock, R., Manolopoulos, Y., 2001. Feature-based classi-  
1120 fication of time-series data. *International Journal of Computer Research*  
1121 10, 49–61.
- 1122 Ngiam, J., Khosla, A., Kim, M., Nam, J., Lee, H., Ng, A.Y., 2011. Multi-  
1123 modal deep learning, in: *In Proceedings of the Twenty-Eighth International*  
1124 *Conference on Machine Learning*.
- 1125 Osuna, G.R., Nagle, T.H., Kermani, B., Schiffman, S.S., 2003. *HandBook of*  
1126 *Machine Olfaction, electronic nose technology*. Wiley-Vch Verlag GmbH &  
1127 Co. KGaA. chapter Signal Conditioning and Preprocessing. pp. 105–132.
- 1128 Pan, W., Torresani, L., 2009. Unsupervised hierarchical modeling of locomo-  
1129 tion styles, in: *Proceedings of the 26th Annual International Conference*  
1130 *on Machine Learning*, pp. 785–792.
- 1131 Parris, E., Carey, M., 1996. Language independent gender identification, in:  
1132 *Acoustics, Speech, and Signal Processing, 1996. ICASSP-96. Conference*  
1133 *Proceedings., 1996 IEEE International Conference on*, pp. 685–688 vol. 2.
- 1134 Pascanu, R., Mikolov, T., Bengio, Y., 2012. Understanding the exploding gra-  
1135 dient problem. *Computing Research Repository (CoRR)* abs/1211.5063.
- 1136 Rabiner, L., Juang, B., 1986. An introduction to hidden markov models.  
1137 *IEEE ASSP Magazine* 3(1), 4–16.
- 1138 Raina, R., Battle, A., Lee, H., Packer, B., Ng, A.Y., 2007. Self-taught  
1139 learning: Transfer learning from unlabeled data, in: *Proceedings of the*  
1140 *Twenty-fourth International Conference on Machine Learning*.

- 1141 Ranzato, M., Hinton, G., 2010. Modeling pixel means and covariances  
1142 using factorized third-order boltzmann machines, in: in Proc. of Computer  
1143 Vision and Pattern Recognition Conference (CVPR 2010).
- 1144 Ranzato, M., Krizhevsky, A., Hinton, G., 2010. Factored 3-way restricted  
1145 boltzmann machines for modeling natural images, in: in Proceedings of  
1146 the International Conference on Artificial Intelligence and Statistics.
- 1147 Ranzato, M., Poultney, C., Chopra, S., LeCun, Y., 2006. Efficient learning of  
1148 sparse representations with an energy-based model, in: et al., J.P. (Ed.),  
1149 Advances in Neural Information Processing Systems (NIPS 2006), MIT  
1150 Press.
- 1151 Saxe, A., Koh, P., Chen, Z., Bhand, M., Suresh, B., Ng, A.Y., 2011. On  
1152 random weights and unsupervised feature learning, in: In Proceedings of  
1153 the Twenty-Eighth International Conference on Machine Learning.
- 1154 Schoerhuber, C., Klapuri, A., 2010. Constant-q transform toolbox for music  
1155 processing, in: 7th Sound and Music Computing Conference.
- 1156 Smith, E., Lewicki, M.S., 2005. Learning efficient auditory codes using spikes  
1157 predicts cochlear filters, in: In Advances in Neural Information Processing  
1158 Systems, MIT Press.
- 1159 Stavens, D., Thrun, S., 2010. Unsupervised learning of invariant features  
1160 using video, in: Computer Vision and Pattern Recognition (CVPR), 2010  
1161 IEEE Conference on, pp. 1649–1656.
- 1162 Sugiyama, M., Sawai, H., Waibel, A., 1991. Review of tdnn (time delay neural

1163 network) architectures for speech recognition, in: Circuits and Systems,  
1164 1991., IEEE International Sympoisum on, pp. 582–585 vol.1.

1165 Sutskever, I., 2012. Training Recurrent Neural Networks. Ph.D. thesis. Uni-  
1166 versity of Toronto.

1167 Sutskever, I., Hinton, G., 2006. Learning multilevel distributed represen-  
1168 tations for high-dimensional sequences. Technical Report. University of  
1169 Toronto.

1170 Sutskever, I., Hinton, G.E., Taylor, G.W., 2008. The recurrent temporal  
1171 restricted boltzmann machine, in: Advances in Neural Information Pro-  
1172 cessing Systems, pp. 1601–1608.

1173 Taylor, G., Fergus, R., LeCun, Y., Bregler, C., 2010. Convolutional learning  
1174 of spatio-temporal features, in: Proc. European Conference on Computer  
1175 Vision (ECCV’10).

1176 Taylor, G., Hinton, G., 2009. Factored conditional restricted boltzmann  
1177 machines for modeling motion style, in: Proc. of the 26th International  
1178 Conference on Machine Learning (ICML).

1179 Taylor, G., Hinton, G.E., Roweis, S., 2007. Modeling human motion using  
1180 binary latent variables, in: Advances in Neural Information Processing  
1181 Systems.

1182 Taylor, G.W., 2009. Composable, distributed-state models for high-  
1183 dimensional time series. Ph.D. thesis. Departmet of Computer Science  
1184 University of Toronto.

- 1185 Trincavelli, M., Coradeschi, S., Loutfi, A., Söderquist, B., Thunberg, P.,  
1186 2010. Direct identification of bacteria in blood culture samples using an  
1187 electronic nose. *IEEE Trans Biomedical Engineering* 57, 2884–2890.
- 1188 Tsai, C.F., Hsiao, Y.C., 2010. Combining multiple feature selection meth-  
1189 ods for stock prediction: Union, intersection, and multi-intersection ap-  
1190 proaches. *Decision Support Systems* 50, 258 – 269.
- 1191 Tucker, C., 1999. Self-organizing maps for time series analysis of electromyo-  
1192 graphic data, in: *Neural Networks, 1999. IJCNN '99. International Joint*  
1193 *Conference on*, pp. 3577–3580.
- 1194 Vembu, S., Vergara, A., Muezzinoglu, M.K., Huerta, R., 2012. On time  
1195 series features and kernels for machine olfaction. *Sensors and Actuators*  
1196 *B: Chemical* 174, 535–546.
- 1197 Vito, S.D., Castaldo, A., Loffredo, F., Massera, E., Polichetti, T., Nasti, I.,  
1198 Vacca, P., Quercia, L., Francia, G.D., 2007. Gas concentration estimation  
1199 in ternary mixtures with room temperature operating sensor array using  
1200 tapped delay architectures. *Sensors and Actuators B: Chemical* 124, 309  
1201 – 316.
- 1202 Waibel, A., Hanazawa, T., Hinton, G., Shikano, K., Lang, K., 1989. Phoneme  
1203 recognition using time-delay neural networks. *IEEE Trans. Acoust.,*  
1204 *Speech, Signal Processing* 37, 328–339.
- 1205 Wang, D., Shang, Y., 2013. Modeling physiological data with deep belief  
1206 networks. *International Journal of Information and Education Technology*  
1207 3.

- 1208 Wang, J.M., Fleet, D.J., Hertzmann, A., 2007. Multi-factor gaussian pro-  
1209 cess models for style-content separation, in: International Conference of  
1210 Machine Learning (ICML), pp. 975–10982.
- 1211 Wiskott, L., Sejnowski, T.J., 2002. Slow feature analysis: Unsupervised  
1212 learning of invariances. *Neural computation* 14, 715–770.
- 1213 Wulsin, D., Gupta, J., Mani, R., Blanco, J., Litt, B., 2011. Modeling  
1214 electroencephalography waveforms with semi-supervised deep belief nets:  
1215 faster classification and anomaly measurement. *Journal of Neural Engi-  
1216 neering* 8, 1741 – 2552.
- 1217 Yamazaki, A., Ludermir, T., De Souto, M.C.P., 2001. Classification of vin-  
1218 tages of wine by artificial nose using time delay neural networks. *Electron-  
1219 ics Letters* 37, 1466–1467.
- 1220 Yang, Q., Wu, X., 2006. 10 challenging problems in data mining research.  
1221 *International Journal of Information Technology & Decision Making* 05,  
1222 597–604.
- 1223 Zampolli, S., Elmi, I., Ahmed, F., Passini, M., Cardinali, G., Nicoletti, S.,  
1224 Dori, L., 2004. An electronic nose based on solid state sensor arrays for  
1225 low-cost indoor air quality monitoring applications. *Sensors and Actuators  
1226 B: Chemical* 101, 39–46.
- 1227 Zhang, H., Balaban, M.O., Principe, J.C., 2003. Improving pattern recogni-  
1228 tion of electronic nose data with time-delay neural networks. *Sensors and  
1229 Actuators B: Chemical* 96, 385–389.

- 1230 Zhu, X., Wang, H., Xu, L., Li, H., 2008. Predicting stock index increments  
1231 by neural networks: The role of trading volume under different horizons.  
1232 Expert Systems with Applications 34, 3043 – 3054.
- 1233 Zou, W.Y., Ng, A.Y., Yu, K., 2011. Unsupervised learning of visual in-  
1234 variance with temporal coherence, in: In NIPS 2011 Workshop on Deep  
1235 Learning and Unsupervised Feature Learning.