

To gain some intuition for how the BPTT algorithm behaves, we provide an example of how to compute gradients by BPTT for the RNN equations above (Eq. 10.8 and Eq. 10.12). The nodes of our computational graph include the parameters  $\mathbf{U}$ ,  $\mathbf{V}$ ,  $\mathbf{W}$ ,  $\mathbf{b}$  and  $\mathbf{c}$  as well as the sequence of nodes indexed by  $t$  for  $\mathbf{x}^{(t)}$ ,  $\mathbf{h}^{(t)}$ ,  $\mathbf{o}^{(t)}$  and  $L^{(t)}$ . For each node  $\mathbf{N}$  we need to compute the gradient  $\nabla_{\mathbf{N}} L$  recursively, based on the gradient computed at nodes that follow it in the graph. We start the recursion with the nodes immediately preceding the final loss

$$\frac{\partial L}{\partial L^{(t)}} = 1. \quad (10.17)$$

In this derivation we assume that the outputs  $\mathbf{o}^{(t)}$  are used as the argument to the softmax function to obtain the vector  $\hat{\mathbf{y}}$  of probabilities over the output. We also assume that the loss is the negative log-likelihood of the true target  $y^{(t)}$  given the input so far. The gradient  $\nabla_{\mathbf{o}^{(t)}} L$  on the outputs at time step  $t$ , for all  $i, t$ , is as follows:

$$(\nabla_{\mathbf{o}^{(t)}} L)_i = \frac{\partial L}{\partial o_i^{(t)}} = \frac{\partial L}{\partial L^{(t)}} \frac{\partial L^{(t)}}{\partial o_i^{(t)}} = \hat{y}_i^{(t)} - \mathbf{1}_{i,y^{(t)}}. \quad (10.18)$$

We work our way backwards, starting from the end of the sequence. At the final time step  $\tau$ ,  $\mathbf{h}^{(\tau)}$  only has  $\mathbf{o}^{(\tau)}$  as a descendent, so its gradient is simple:

$$\nabla_{\mathbf{h}^{(\tau)}} L = (\nabla_{\mathbf{o}^{(\tau)}} L) \frac{\partial \mathbf{o}^{(\tau)}}{\partial \mathbf{h}^{(\tau)}} = (\nabla_{\mathbf{o}^{(\tau)}} L) \mathbf{V}. \quad (10.19)$$

We can then iterate backwards in time to back-propagate gradients through time, from  $t = \tau - 1$  down to  $t = 1$ , noting that  $\mathbf{h}^{(t)}$  (for  $t < \tau$ ) has as descendants both  $\mathbf{o}^{(t)}$  and  $\mathbf{h}^{(t+1)}$ . Its gradient is thus given by

$$\nabla_{\mathbf{h}^{(t)}} L = (\nabla_{\mathbf{h}^{(t+1)}} L) \frac{\partial \mathbf{h}^{(t+1)}}{\partial \mathbf{h}^{(t)}} + (\nabla_{\mathbf{o}^{(t)}} L) \frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{h}^{(t)}} \quad (10.20)$$

$$= (\nabla_{\mathbf{h}^{(t+1)}} L) \text{diag} \left( 1 - (\mathbf{h}^{(t+1)})^2 \right) \mathbf{W} + (\nabla_{\mathbf{o}^{(t)}} L) \mathbf{V} \quad (10.21)$$

where  $\text{diag} \left( 1 - (\mathbf{h}^{(t+1)})^2 \right)$  indicates the diagonal matrix containing the elements  $1 - (h_i^{(t+1)})^2$ . This is the Jacobian of the hyperbolic tangent associated with the hidden unit  $i$  at time  $t + 1$ .

Once the gradients on the internal nodes of the computational graph are obtained, we can obtain the gradients on the parameter nodes, which have descendants at all the time steps:

$$\nabla_{\mathbf{c}} L = \sum_t (\nabla_{\mathbf{o}^{(t)}} L) \frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{c}} = \sum_t \nabla_{\mathbf{o}^{(t)}} L$$

$$\begin{aligned}
 \nabla_{\mathbf{b}} L &= \sum_t (\nabla_{\mathbf{h}^{(t)}} L) \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{b}} = \sum_t (\nabla_{\mathbf{h}^{(t)}} L) \operatorname{diag} \left( 1 - (\mathbf{h}^{(t)})^2 \right) \\
 \nabla_{\mathbf{V}} L &= \sum_t (\nabla_{\mathbf{o}^{(t)}} L) \frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{V}} = \sum_t (\nabla_{\mathbf{o}^{(t)}} L) \mathbf{h}^{(t)\top} \\
 \nabla_{\mathbf{W}} L &= \sum_t (\nabla_{\mathbf{h}^{(t)}} L) \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{W}} = \sum_t (\nabla_{\mathbf{h}^{(t)}} L) \operatorname{diag} \left( 1 - (\mathbf{h}^{(t)})^2 \right) \mathbf{h}^{(t-1)\top} \\
 \nabla_{\mathbf{U}} L &= \sum_t (\nabla_{\mathbf{h}^{(t)}} L) \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{U}} = \sum_t (\nabla_{\mathbf{h}^{(t)}} L) \operatorname{diag} \left( 1 - (\mathbf{h}^{(t)})^2 \right) \mathbf{x}^{(t)\top}
 \end{aligned}$$

We do not need to compute the gradient with respect to  $\mathbf{x}^{(t)}$  for training because it does not have any parameters as ancestors in the computational graph defining the loss.

We are abusing notation somewhat in the above equations. We correctly use  $\nabla_{\mathbf{h}^{(t)}} L$  to indicate the full influence of  $\mathbf{h}^{(t)}$  through all paths from  $\mathbf{h}^{(t)}$  to  $L$ . This is in contrast to our usage of  $\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{W}}$  or  $\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{b}}$ , which we use here in an unconventional manner. By  $\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{W}}$  we refer to the effect of  $\mathbf{W}$  on  $\mathbf{h}^{(t)}$  only via the use of  $\mathbf{W}$  at time step  $t$ . This is not standard calculus notation, because the standard definition of the Jacobian would actually include the complete influence of  $\mathbf{W}$  on  $\mathbf{h}^{(t)}$  via its use in all of the preceding time steps to produce  $\mathbf{h}^{(t-1)}$ . What we refer to here is in fact the `bprop` method of Sec. 6.5.6, that computes the contribution of a single edge in the computational graph to the gradient.

### 10.2.3 Recurrent Networks as Directed Graphical Models

In the example recurrent network we have developed so far, the losses  $L^{(t)}$  were cross-entropies between training targets  $\mathbf{y}^{(t)}$  and outputs  $\mathbf{o}^{(t)}$ . As with a feedforward network, it is in principle possible to use almost any loss with a recurrent network. The loss should be chosen based on the task. As with a feedforward network, we usually wish to interpret the output of the RNN as a probability distribution, and we usually use the cross-entropy associated with that distribution to define the loss. Mean squared error is the cross-entropy loss associated with an output distribution that is a unit Gaussian, for example, just as with a feedforward network.

When we use a predictive log-likelihood training objective, such as Eq. 10.12, we train the RNN to estimate the conditional distribution of the next sequence element  $\mathbf{y}^{(t)}$  given the past inputs. This may mean that we maximize the log-likelihood

$$\log p(\mathbf{y}^{(t)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}), \quad (10.22)$$

or, if the model includes connections from the output at one time step to the next time step,

$$\log p(\mathbf{y}^{(t)} \mid \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}, \mathbf{y}^{(1)}, \dots, \mathbf{y}^{(t-1)}). \quad (10.23)$$

Decomposing the joint probability over the sequence of  $\mathbf{y}$  values as a series of one-step probabilistic predictions is one way to capture the full joint distribution across the whole sequence. When we do not feed past  $\mathbf{y}$  values as inputs that condition the next step prediction, the directed graphical model contains no edges from any  $\mathbf{y}^{(i)}$  in the past to the current  $\mathbf{y}^{(t)}$ . In this case, the outputs  $\mathbf{y}$  are conditionally independent given the sequence of  $\mathbf{x}$  values. When we do feed the actual  $\mathbf{y}$  values (not their prediction, but the actual observed or generated values) back into the network, the directed graphical model contains edges from all  $\mathbf{y}^{(i)}$  values in the past to the current  $\mathbf{y}^{(t)}$  value.

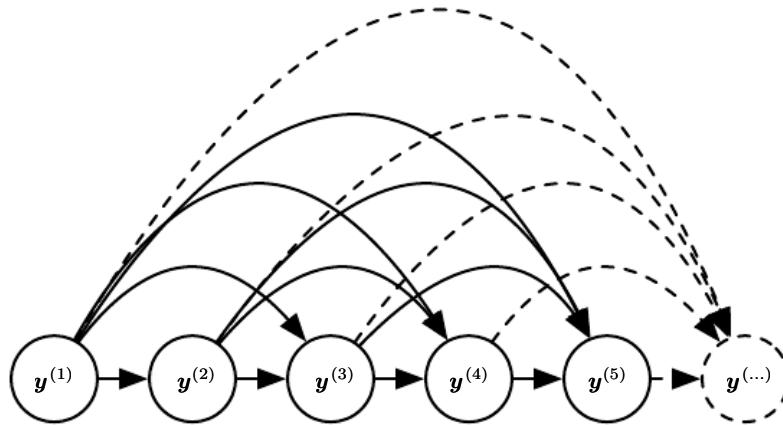


Figure 10.7: Fully connected graphical model for a sequence  $y^{(1)}, y^{(2)}, \dots, y^{(t)}, \dots$ : every past observation  $y^{(i)}$  may influence the conditional distribution of some  $y^{(t)}$  (for  $t > i$ ), given the previous values. Parametrizing the graphical model directly according to this graph (as in Eq. 10.6) might be very inefficient, with an ever growing number of inputs and parameters for each element of the sequence. RNNs obtain the same full connectivity but efficient parametrization, as illustrated in Fig. 10.8.

As a simple example, let us consider the case where the RNN models only a sequence of scalar random variables  $\mathbb{Y} = \{y^{(1)}, \dots, y^{(\tau)}\}$ , with no additional inputs  $\mathbf{x}$ . The input at time step  $t$  is simply the output at time step  $t - 1$ . The RNN then defines a directed graphical model over the  $y$  variables. We parametrize the joint distribution of these observations using the chain rule (Eq. 3.6) for conditional probabilities:

$$P(\mathbb{Y}) = P(\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)}) = \prod_{t=1}^{\tau} P(y^{(t)} \mid \mathbf{y}^{(t-1)}, \mathbf{y}^{(t-2)}, \dots, \mathbf{y}^{(1)}) \quad (10.24)$$

where the right-hand side of the bar is empty for  $t = 1$ , of course. Hence the negative log-likelihood of a set of values  $\{y^{(1)}, \dots, y^{(\tau)}\}$  according to such a model is

$$L = \sum_t L^{(t)} \quad (10.25)$$

where

$$L^{(t)} = -\log P(y^{(t)} = y^{(t)} | y^{(t-1)}, y^{(t-2)}, \dots, y^{(1)}). \quad (10.26)$$

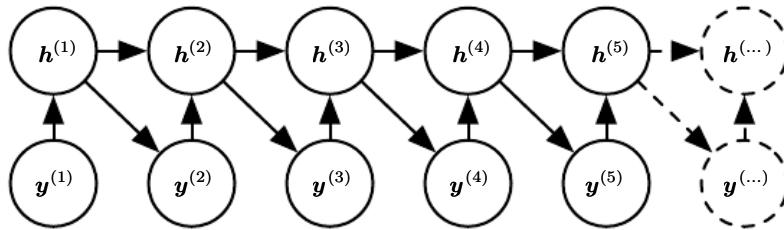


Figure 10.8: Introducing the state variable in the graphical model of the RNN, even though it is a deterministic function of its inputs, helps to see how we can obtain a very efficient parametrization, based on Eq. 10.5. Every stage in the sequence (for  $h^{(t)}$  and  $y^{(t)}$ ) involves the same structure (the same number of inputs for each node) and can share the same parameters with the other stages.

The edges in a graphical model indicate which variables depend directly on other variables. Many graphical models aim to achieve statistical and computational efficiency by omitting edges that do not correspond to strong interactions. For example, it is common to make the Markov assumption that the graphical model should only contain edges from  $\{y^{(t-k)}, \dots, y^{(t-1)}\}$  to  $y^{(t)}$ , rather than containing edges from the entire past history. However, in some cases, we believe that all past inputs should have an influence on the next element of the sequence. RNNs are useful when we believe that the distribution over  $y^{(t)}$  may depend on a value of  $y^{(i)}$  from the distant past in a way that is not captured by the effect of  $y^{(i)}$  on  $y^{(t-1)}$ .

One way to interpret an RNN as a graphical model is to view the RNN as defining a graphical model whose structure is the complete graph, able to represent direct dependencies between any pair of  $y$  values. The graphical model over the  $y$  values with the complete graph structure is shown in Fig. 10.7. The complete graph interpretation of the RNN is based on ignoring the hidden units  $h^{(t)}$  by marginalizing them out of the model.

It is more interesting to consider the graphical model structure of RNNs that results from regarding the hidden units  $h^{(t)}$  as random variables.<sup>1</sup> Including the

<sup>1</sup>The conditional distribution over these variables given their parents is deterministic. This is

hidden units in the graphical model reveals that the RNN provides a very efficient parametrization of the joint distribution over the observations. Suppose that we represented an arbitrary joint distribution over discrete values with a tabular representation—an array containing a separate entry for each possible assignment of values, with the value of that entry giving the probability of that assignment occurring. If  $y$  can take on  $k$  different values, the tabular representation would have  $O(k^\tau)$  parameters. By comparison, due to parameter sharing, the number of parameters in the RNN is  $O(1)$  as a function of sequence length. The number of parameters in the RNN may be adjusted to control model capacity but is not forced to scale with sequence length. Eq. 10.5 shows that the RNN parametrizes long-term relationships between variables efficiently, using recurrent applications of the same function  $f$  and same parameters  $\theta$  at each time step. Fig. 10.8 illustrates the graphical model interpretation. Incorporating the  $\mathbf{h}^{(t)}$  nodes in the graphical model decouples the past and the future, acting as an intermediate quantity between them. A variable  $y^{(i)}$  in the distant past may influence a variable  $y^{(t)}$  via its effect on  $\mathbf{h}$ . The structure of this graph shows that the model can be efficiently parametrized by using the same conditional probability distributions at each time step, and that when the variables are all observed, the probability of the joint assignment of all variables can be evaluated efficiently.

Even with the efficient parametrization of the graphical model, some operations remain computationally challenging. For example, it is difficult to predict missing values in the middle of the sequence.

The price recurrent networks pay for their reduced number of parameters is that **optimizing** the parameters may be difficult.

The parameter sharing used in recurrent networks relies on the assumption that the same parameters can be used for different time steps. Equivalently, the assumption is that the conditional probability distribution over the variables at time  $t + 1$  given the variables at time  $t$  is *stationary*, meaning that the relationship between the previous time step and the next time step does not depend on  $t$ . In principle, it would be possible to use  $t$  as an extra input at each time step and let the learner discover any time-dependence while sharing as much as it can between different time steps. This would already be much better than using a different conditional probability distribution for each  $t$ , but the network would then have to extrapolate when faced with new values of  $t$ .

To complete our view of an RNN as a graphical model, we must describe how to draw samples from the model. The main operation that we need to perform is

---

perfectly legitimate, though it is somewhat rare to design a graphical model with such deterministic hidden units.

simply to sample from the conditional distribution at each time step. However, there is one additional complication. The RNN must have some mechanism for determining the length of the sequence. This can be achieved in various ways.

In the case when the output is a symbol taken from a vocabulary, one can add a special symbol corresponding to the end of a sequence (Schmidhuber, 2012). When that symbol is generated, the sampling process stops. In the training set, we insert this symbol as an extra member of the sequence, immediately after  $\mathbf{x}^{(\tau)}$  in each training example.

Another option is to introduce an extra Bernoulli output to the model that represents the decision to either continue generation or halt generation at each time step. This approach is more general than the approach of adding an extra symbol to the vocabulary, because it may be applied to any RNN, rather than only RNNs that output a sequence of symbols. For example, it may be applied to an RNN that emits a sequence of real numbers. The new output unit is usually a sigmoid unit trained with the cross-entropy loss. In this approach the sigmoid is trained to maximize the log-probability of the correct prediction as to whether the sequence ends or continues at each time step.

Another way to determine the sequence length  $\tau$  is to add an extra output to the model that predicts the integer  $\tau$  itself. The model can sample a value of  $\tau$  and then sample  $\tau$  steps worth of data. This approach requires adding an extra input to the recurrent update at each time step so that the recurrent update is aware of whether it is near the end of the generated sequence. This extra input can either consist of the value of  $\tau$  or can consist of  $\tau - t$ , the number of remaining time steps. Without this extra input, the RNN might generate sequences that end abruptly, such as a sentence that ends before it is complete. This approach is based on the decomposition

$$P(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)}) = P(\tau)P(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)} \mid \tau). \quad (10.27)$$

The strategy of predicting  $\tau$  directly is used for example by Goodfellow *et al.* (2014d).

#### 10.2.4 Modeling Sequences Conditioned on Context with RNNs

In the previous section we described how an RNN could correspond to a directed graphical model over a sequence of random variables  $y^{(t)}$  with no inputs  $\mathbf{x}$ . Of course, our development of RNNs as in Eq. 10.8 included a sequence of inputs  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(\tau)}$ . In general, RNNs allow the extension of the graphical model view to represent not only a joint distribution over the  $y$  variables but also a

conditional distribution over  $y$  given  $\mathbf{x}$ . As discussed in the context of feedforward networks in Sec. 6.2.1.1, any model representing a variable  $P(\mathbf{y}; \boldsymbol{\theta})$  can be reinterpreted as a model representing a conditional distribution  $P(\mathbf{y}|\boldsymbol{\omega})$  with  $\boldsymbol{\omega} = \boldsymbol{\theta}$ . We can extend such a model to represent a distribution  $P(\mathbf{y} | \mathbf{x})$  by using the same  $P(\mathbf{y} | \boldsymbol{\omega})$  as before, but making  $\boldsymbol{\omega}$  a function of  $\mathbf{x}$ . In the case of an RNN, this can be achieved in different ways. We review here the most common and obvious choices.

Previously, we have discussed RNNs that take a sequence of vectors  $\mathbf{x}^{(t)}$  for  $t = 1, \dots, \tau$  as input. Another option is to take only a single vector  $\mathbf{x}$  as input. When  $\mathbf{x}$  is a fixed-size vector, we can simply make it an extra input of the RNN that generates the  $\mathbf{y}$  sequence. Some common ways of providing an extra input to an RNN are:

1. as an extra input at each time step, or
2. as the initial state  $\mathbf{h}^{(0)}$ , or
3. both.

The first and most common approach is illustrated in Fig. 10.9. The interaction between the input  $\mathbf{x}$  and each hidden unit vector  $\mathbf{h}^{(t)}$  is parametrized by a newly introduced weight matrix  $\mathbf{R}$  that was absent from the model of only the sequence of  $y$  values. The same product  $\mathbf{x}^\top \mathbf{R}$  is added as additional input to the hidden units at every time step. We can think of the choice of  $\mathbf{x}$  as determining the value of  $\mathbf{x}^\top \mathbf{R}$  that is effectively a new bias parameter used for each of the hidden units. The weights remain independent of the input. We can think of this model as taking the parameters  $\boldsymbol{\theta}$  of the non-conditional model and turning them into  $\boldsymbol{\omega}$ , where the bias parameters within  $\boldsymbol{\omega}$  are now a function of the input.

Rather than receiving only a single vector  $\mathbf{x}$  as input, the RNN may receive a sequence of vectors  $\mathbf{x}^{(t)}$  as input. The RNN described in Eq. 10.8 corresponds to a conditional distribution  $P(\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)})$  that makes a conditional independence assumption that this distribution factorizes as

$$\prod_t P(\mathbf{y}^{(t)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}). \quad (10.28)$$

To remove the conditional independence assumption, we can add connections from the output at time  $t$  to the hidden unit at time  $t+1$ , as shown in Fig. 10.10. The model can then represent arbitrary probability distributions over the  $\mathbf{y}$  sequence. This kind of model representing a distribution over a sequence given another sequence still has one restriction, which is that the length of both sequences must be the same. We describe how to remove this restriction in Sec. 10.4.

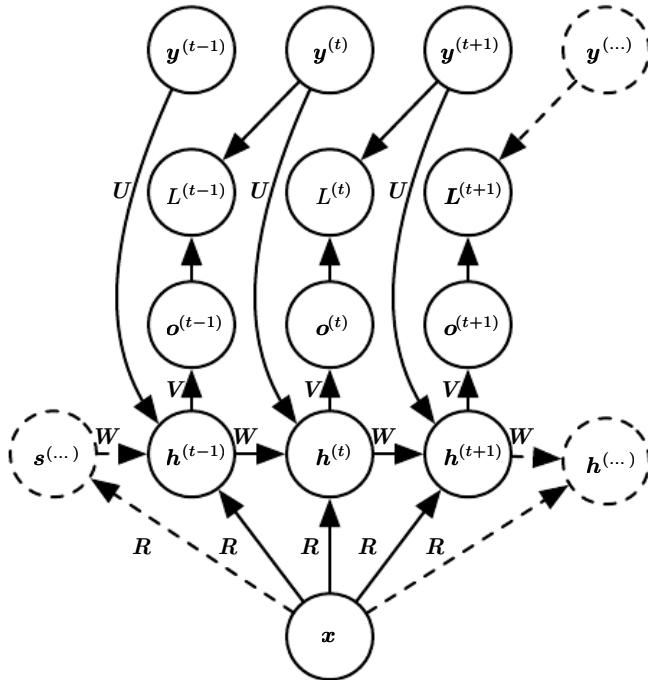


Figure 10.9: An RNN that maps a fixed-length vector  $x$  into a distribution over sequences  $\mathbf{Y}$ . This RNN is appropriate for tasks such as image captioning, where a single image is used as input to a model that then produces a sequence of words describing the image. Each element  $y^{(t)}$  of the observed output sequence serves both as input (for the current time step) and, during training, as target (for the previous time step).

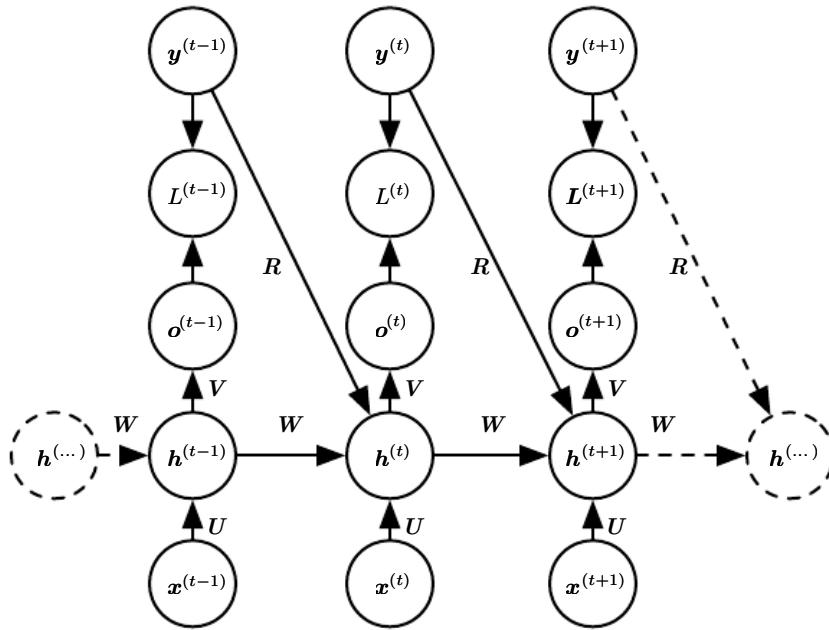


Figure 10.10: A conditional recurrent neural network mapping a variable-length sequence of  $\mathbf{x}$  values into a distribution over sequences of  $\mathbf{y}$  values of the same length. Compared to Fig. 10.3, this RNN contains connections from the previous output to the current state. These connections allow this RNN to model an arbitrary distribution over sequences of  $\mathbf{y}$  given sequences of  $\mathbf{x}$  of the same length. The RNN of Fig. 10.3 is only able to represent distributions in which the  $\mathbf{y}$  values are conditionally independent from each other given the  $\mathbf{x}$  values.

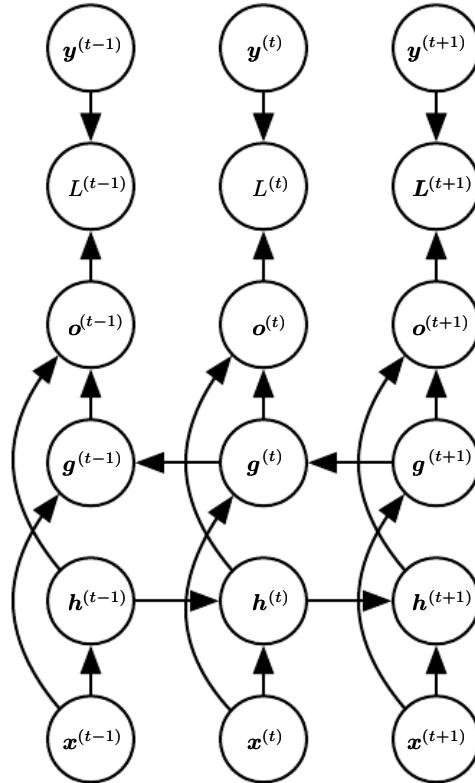


Figure 10.11: Computation of a typical bidirectional recurrent neural network, meant to learn to map input sequences  $\mathbf{x}$  to target sequences  $\mathbf{y}$ , with loss  $L^{(t)}$  at each step  $t$ . The  $\mathbf{h}$  recurrence propagates information forward in time (towards the right) while the  $\mathbf{g}$  recurrence propagates information backward in time (towards the left). Thus at each point  $t$ , the output units  $\mathbf{o}^{(t)}$  can benefit from a relevant summary of the past in its  $\mathbf{h}^{(t)}$  input and from a relevant summary of the future in its  $\mathbf{g}^{(t)}$  input.

## 10.3 Bidirectional RNNs

All of the recurrent networks we have considered up to now have a “causal” structure, meaning that the state at time  $t$  only captures information from the past,  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t-1)}$ , and the present input  $\mathbf{x}^{(t)}$ . Some of the models we have discussed also allow information from past  $\mathbf{y}$  values to affect the current state when the  $\mathbf{y}$  values are available.

However, in many applications we want to output a prediction of  $\mathbf{y}^{(t)}$  which may depend on **the whole input sequence**. For example, in speech recognition, the correct interpretation of the current sound as a phoneme may depend on the next few phonemes because of co-articulation and potentially may even depend on the next few words because of the linguistic dependencies between nearby words: if there are two interpretations of the current word that are both acoustically plausible, we may have to look far into the future (and the past) to disambiguate them. This is also true of handwriting recognition and many other sequence-to-sequence learning tasks, described in the next section.

Bidirectional recurrent neural networks (or bidirectional RNNs) were invented to address that need (Schuster and Paliwal, 1997). They have been extremely successful (Graves, 2012) in applications where that need arises, such as handwriting recognition (Graves *et al.*, 2008; Graves and Schmidhuber, 2009), speech recognition (Graves and Schmidhuber, 2005; Graves *et al.*, 2013) and bioinformatics (Baldi *et al.*, 1999).

As the name suggests, bidirectional RNNs combine an RNN that moves forward through time beginning from the start of the sequence with another RNN that moves backward through time beginning from the end of the sequence. Fig. 10.11 illustrates the typical bidirectional RNN, with  $\mathbf{h}^{(t)}$  standing for the state of the sub-RNN that moves forward through time and  $\mathbf{g}^{(t)}$  standing for the state of the sub-RNN that moves backward through time. This allows the output units  $\mathbf{o}^{(t)}$  to compute a representation that depends on **both the past and the future** but is most sensitive to the input values around time  $t$ , without having to specify a fixed-size window around  $t$  (as one would have to do with a feedforward network, a convolutional network, or a regular RNN with a fixed-size look-ahead buffer).

This idea can be naturally extended to 2-dimensional input, such as images, by having **four** RNNs, each one going in one of the four directions: up, down, left, right. At each point  $(i, j)$  of a 2-D grid, an output  $O_{i,j}$  could then compute a representation that would capture mostly local information but could also depend on long-range inputs, if the RNN is able to learn to carry that information. Compared to a convolutional network, RNNs applied to images are typically more

expensive but allow for long-range lateral interactions between features in the same feature map (Visin *et al.*, 2015; Kalchbrenner *et al.*, 2015). Indeed, the forward propagation equations for such RNNs may be written in a form that shows they use a convolution that computes the bottom-up input to each layer, prior to the recurrent propagation across the feature map that incorporates the lateral interactions.

## 10.4 Encoder-Decoder Sequence-to-Sequence Architectures

We have seen in Fig. 10.5 how an RNN can map an input sequence to a fixed-size vector. We have seen in Fig. 10.9 how an RNN can map a fixed-size vector to a sequence. We have seen in Fig. 10.3, Fig. 10.4, Fig. 10.10 and Fig. 10.11 how an RNN can map an input sequence to an output sequence of the same length.

Here we discuss how an RNN can be trained to map an input sequence to an output sequence which is not necessarily of the same length. This comes up in many applications, such as speech recognition, machine translation or question answering, where the input and output sequences in the training set are generally not of the same length (although their lengths might be related).

We often call the input to the RNN the “context.” We want to produce a representation of this context,  $C$ . The context  $C$  might be a vector or sequence of vectors that summarize the input sequence  $\mathbf{X} = (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n_x)})$ .

The simplest RNN architecture for mapping a variable-length sequence to another variable-length sequence was first proposed by Cho *et al.* (2014a) and shortly after by Sutskever *et al.* (2014), who independently developed that architecture and were the first to obtain state-of-the-art translation using this approach. The former system is based on scoring proposals generated by another machine translation system, while the latter uses a standalone recurrent network to generate the translations. These authors respectively called this architecture, illustrated in Fig. 10.12, the encoder-decoder or sequence-to-sequence architecture. The idea is very simple: (1) an *encoder* or *reader* or *input* RNN processes the input sequence. The encoder emits the context  $C$ , usually as a simple function of its final hidden state. (2) a *decoder* or *writer* or *output* RNN is conditioned on that fixed-length vector (just like in Fig. 10.9) to generate the output sequence  $\mathbf{Y} = (\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(n_y)})$ . The innovation of this kind of architecture over those presented in earlier sections of this chapter is that the lengths  $n_x$  and  $n_y$  can vary from each other, while previous architectures constrained  $n_x = n_y = \tau$ . In a sequence-to-sequence architecture, the two RNNs

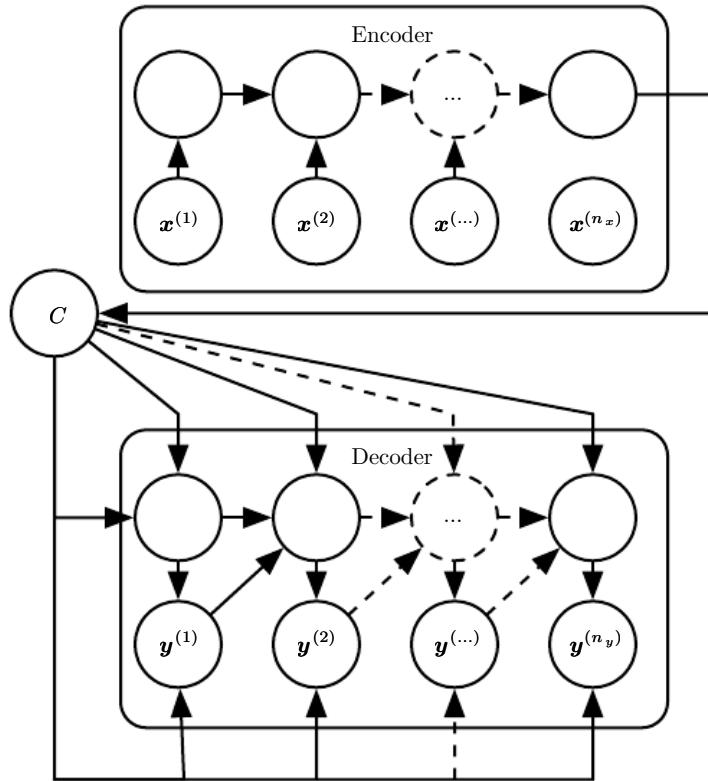


Figure 10.12: Example of an encoder-decoder or sequence-to-sequence RNN architecture, for learning to generate an output sequence  $(\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(n_y)})$  given an input sequence  $(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n_x)})$ . It is composed of an encoder RNN that reads the input sequence and a decoder RNN that generates the output sequence (or computes the probability of a given output sequence). The final hidden state of the encoder RNN is used to compute a generally fixed-size context variable  $C$  which represents a semantic summary of the input sequence and is given as input to the decoder RNN.

are trained jointly to maximize the average of  $\log P(\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(n_y)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n_x)})$  over all the pairs of  $\mathbf{x}$  and  $\mathbf{y}$  sequences in the training set. The last state  $\mathbf{h}_{n_x}$  of the encoder RNN is typically used as a representation  $C$  of the input sequence that is provided as input to the decoder RNN.

If the context  $C$  is a vector, then the decoder RNN is simply a vector-to-sequence RNN as described in Sec. 10.2.4. As we have seen, there are at least two ways for a vector-to-sequence RNN to receive input. The input can be provided as the initial state of the RNN, or the input can be connected to the hidden units at each time step. These two ways can also be combined.

There is no constraint that the encoder must have the same size of hidden layer as the decoder.

One clear limitation of this architecture is when the context  $C$  output by the encoder RNN has a dimension that is too small to properly summarize a long sequence. This phenomenon was observed by Bahdanau *et al.* (2015) in the context of machine translation. They proposed to make  $C$  a variable-length sequence rather than a fixed-size vector. Additionally, they introduced an *attention mechanism* that learns to associate elements of the sequence  $C$  to elements of the output sequence. See Sec. 12.4.5.1 for more details.

## 10.5 Deep Recurrent Networks

The computation in most RNNs can be decomposed into three blocks of parameters and associated transformations:

1. from the input to the hidden state,
2. from the previous hidden state to the next hidden state, and
3. from the hidden state to the output.

With the RNN architecture of Fig. 10.3, each of these three blocks is associated with a single weight matrix. In other words, when the network is unfolded, each of these corresponds to a shallow transformation. By a shallow transformation, we mean a transformation that would be represented by a single layer within a deep MLP. Typically this is a transformation represented by a learned affine transformation followed by a fixed nonlinearity.

Would it be advantageous to introduce depth in each of these operations? Experimental evidence (Graves *et al.*, 2013; Pascanu *et al.*, 2014a) strongly suggests so. The experimental evidence is in agreement with the idea that we need enough

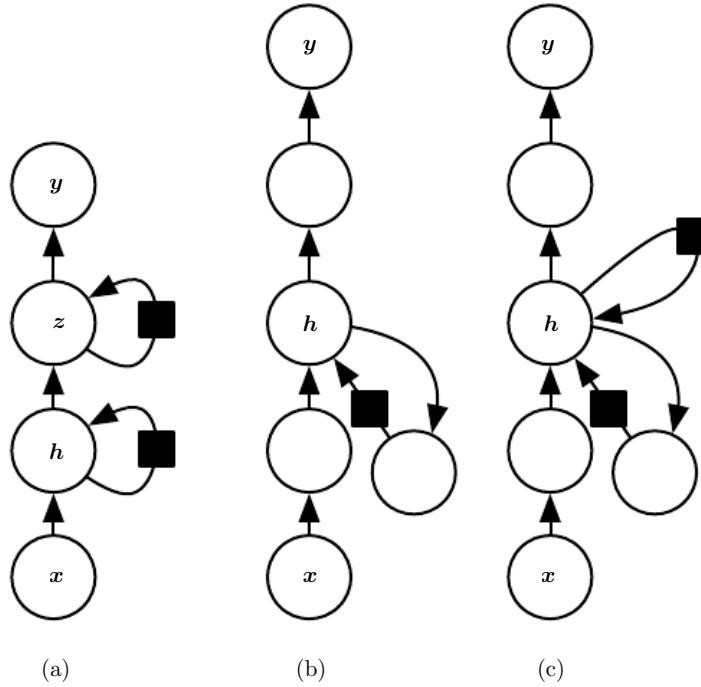


Figure 10.13: A recurrent neural network can be made deep in many ways (Pascanu *et al.*, 2014a). (a) The hidden recurrent state can be broken down into groups organized hierarchically. (b) Deeper computation (e.g., an MLP) can be introduced in the input-to-hidden, hidden-to-hidden and hidden-to-output parts. This may lengthen the shortest path linking different time steps. (c) The path-lengthening effect can be mitigated by introducing skip connections.

depth in order to perform the required mappings. See also Schmidhuber (1992), El Hihi and Bengio (1996), or Jaeger (2007a) for earlier work on deep RNNs.

Graves *et al.* (2013) were the first to show a significant benefit of decomposing the state of an RNN into multiple layers as in Fig. 10.13 (left). We can think of the lower layers in the hierarchy depicted in Fig. 10.13a as playing a role in transforming the raw input into a representation that is more appropriate, at the higher levels of the hidden state. Pascanu *et al.* (2014a) go a step further and propose to have a separate MLP (possibly deep) for each of the three blocks enumerated above, as illustrated in Fig. 10.13b. Considerations of representational capacity suggest to allocate enough capacity in each of these three steps, but doing so by adding depth may hurt learning by making optimization difficult. In general, it is easier to optimize shallower architectures, and adding the extra depth of Fig. 10.13b makes the shortest path from a variable in time step  $t$  to a variable in time step  $t + 1$  become longer. For example, if an MLP with a single hidden layer is used for the state-to-state transition, we have doubled the length of the shortest path between variables in any two different time steps, compared with the ordinary RNN of Fig. 10.3. However, as argued by Pascanu *et al.* (2014a), this can be mitigated by introducing skip connections in the hidden-to-hidden path, as illustrated in Fig. 10.13c.

## 10.6 Recursive Neural Networks

Recursive neural networks<sup>2</sup> represent yet another generalization of recurrent networks, with a different kind of computational graph, which is structured as a deep tree, rather than the chain-like structure of RNNs. The typical computational graph for a recursive network is illustrated in Fig. 10.14. Recursive neural networks were introduced by Pollack (1990) and their potential use for learning to reason was described by Bottou (2011). Recursive networks have been successfully applied to processing **data structures** as input to neural nets (Frasconi *et al.*, 1997, 1998), in natural language processing (Socher *et al.*, 2011a,c, 2013a) as well as in computer vision (Socher *et al.*, 2011b).

One clear advantage of recursive nets over recurrent nets is that for a sequence of the same length  $\tau$ , the depth (measured as the number of compositions of nonlinear operations) can be drastically reduced from  $\tau$  to  $O(\log \tau)$ , which might help deal with long-term dependencies. An open question is how to best structure the tree. One option is to have a tree structure which does not depend on the data,

---

<sup>2</sup>We suggest to not abbreviate “recursive neural network” as “RNN” to avoid confusion with “recurrent neural network.”

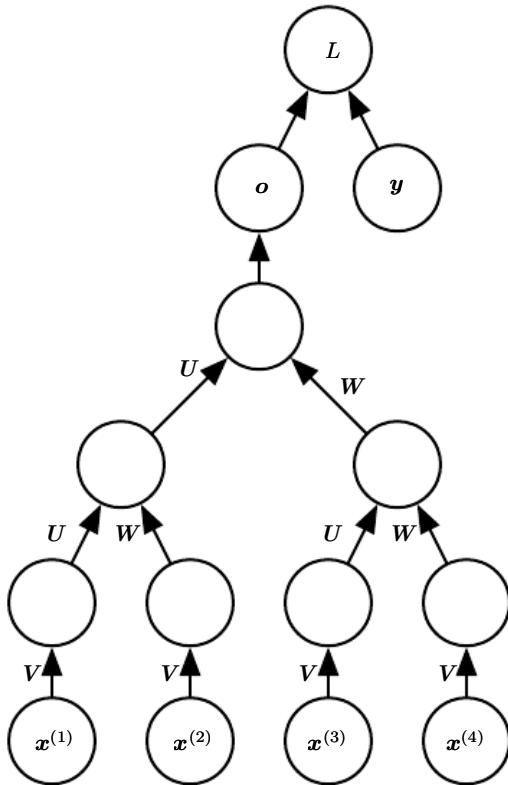


Figure 10.14: A recursive network has a computational graph that generalizes that of the recurrent network from a chain to a tree. A variable-size sequence  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(t)}$  can be mapped to a fixed-size representation (the output  $\mathbf{o}$ ), with a fixed set of parameters (the weight matrices  $\mathbf{U}, \mathbf{V}, \mathbf{W}$ ). The figure illustrates a supervised learning case in which some target  $\mathbf{y}$  is provided which is associated with the whole sequence.

such as a balanced binary tree. In some application domains, external methods can suggest the appropriate tree structure. For example, when processing natural language sentences, the tree structure for the recursive network can be fixed to the structure of the parse tree of the sentence provided by a natural language parser (Socher *et al.*, 2011a, 2013a). Ideally, one would like the learner itself to discover and infer the tree structure that is appropriate for any given input, as suggested by Bottou (2011).

Many variants of the recursive net idea are possible. For example, Frasconi *et al.* (1997) and Frasconi *et al.* (1998) associate the data with a tree structure, and associate the inputs and targets with individual nodes of the tree. The computation performed by each node does not have to be the traditional artificial neuron computation (affine transformation of all inputs followed by a monotone nonlinearity). For example, Socher *et al.* (2013a) propose using tensor operations and bilinear forms, which have previously been found useful to model relationships between concepts (Weston *et al.*, 2010; Bordes *et al.*, 2012) when the concepts are represented by continuous vectors (embeddings).

## 10.7 The Challenge of Long-Term Dependencies

The mathematical challenge of learning long-term dependencies in recurrent networks was introduced in Sec. 8.2.5. The basic problem is that gradients propagated over many stages tend to either vanish (most of the time) or explode (rarely, but with much damage to the optimization). Even if we assume that the parameters are such that the recurrent network is stable (can store memories, with gradients not exploding), the difficulty with long-term dependencies arises from the exponentially smaller weights given to long-term interactions (involving the multiplication of many Jacobians) compared to short-term ones. Many other sources provide a deeper treatment (Hochreiter, 1991; Doya, 1993; Bengio *et al.*, 1994; Pascanu *et al.*, 2013a). In this section, we describe the problem in more detail. The remaining sections describe approaches to overcoming the problem.

Recurrent networks involve the composition of the same function multiple times, once per time step. These compositions can result in extremely nonlinear behavior, as illustrated in Fig. 10.15.

In particular, the function composition employed by recurrent neural networks somewhat resembles matrix multiplication. We can think of the recurrence relation

$$\mathbf{h}^{(t)} = \mathbf{W}^\top \mathbf{h}^{(t-1)} \quad (10.29)$$

as a very simple recurrent neural network lacking a nonlinear activation function,

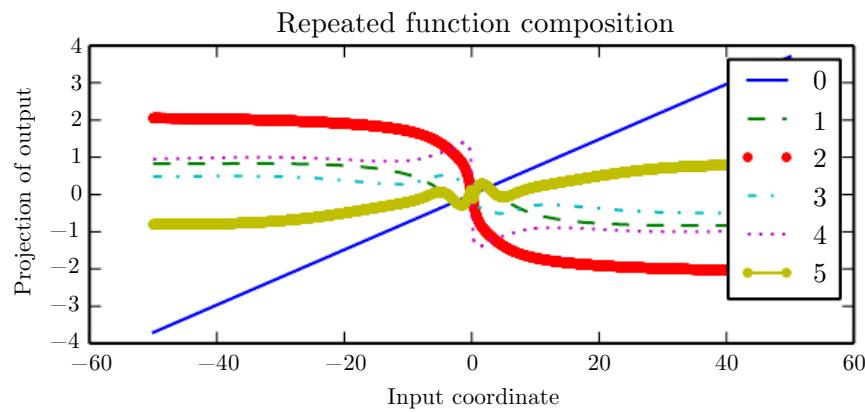


Figure 10.15: When composing many nonlinear functions (like the linear-tanh layer shown here), the result is highly nonlinear, typically with most of the values associated with a tiny derivative, some values with a large derivative, and many alternations between increasing and decreasing. In this plot, we plot a linear projection of a 100-dimensional hidden state down to a single dimension, plotted on the  $y$ -axis. The  $x$ -axis is the coordinate of the initial state along a random direction in the 100-dimensional space. We can thus view this plot as a linear cross-section of a high-dimensional function. The plots show the function after each time step, or equivalently, after each number of times the transition function has been composed.

and lacking inputs  $\mathbf{x}$ . As described in Sec. 8.2.5, this recurrence relation essentially describes the power method. It may be simplified to

$$\mathbf{h}^{(t)} = (\mathbf{W}^t)^\top \mathbf{h}^{(0)}, \quad (10.30)$$

and if  $\mathbf{W}$  admits an eigendecomposition

$$\mathbf{W} = \mathbf{Q}\Lambda\mathbf{Q}^\top, \quad (10.31)$$

the recurrence may be simplified further to

$$\mathbf{h}^{(t)} = \mathbf{Q}^\top \Lambda^t \mathbf{Q} \mathbf{h}^{(0)}. \quad (10.32)$$

The eigenvalues are raised to the power of  $t$  causing eigenvalues with magnitude less than one to decay to zero and eigenvalues with magnitude greater than one to explode. Any component of  $\mathbf{h}^{(0)}$  that is not aligned with the largest eigenvector will eventually be discarded.

This problem is particular to recurrent networks. In the scalar case, imagine multiplying a weight  $w$  by itself many times. The product  $w^t$  will either vanish or explode depending on the magnitude of  $w$ . However, if we make a non-recurrent network that has a different weight  $w^{(t)}$  at each time step, the situation is different. If the initial state is given by 1, then the state at time  $t$  is given by  $\prod_t w^{(t)}$ . Suppose that the  $w^{(t)}$  values are generated randomly, independently from one another, with zero mean and variance  $v$ . The variance of the product is  $O(v^n)$ . To obtain some desired variance  $v^*$  we may choose the individual weights with variance  $v = \sqrt[n]{v^*}$ . Very deep feedforward networks with carefully chosen scaling can thus avoid the vanishing and exploding gradient problem, as argued by [Sussillo \(2014\)](#).

The vanishing and exploding gradient problem for RNNs was independently discovered by separate researchers ([Hochreiter, 1991](#); [Bengio et al., 1993, 1994](#)). One may hope that the problem can be avoided simply by staying in a region of parameter space where the gradients do not vanish or explode. Unfortunately, in order to store memories in a way that is robust to small perturbations, the RNN must enter a region of parameter space where gradients vanish ([Bengio et al., 1993, 1994](#)). Specifically, whenever the model is able to represent long term dependencies, the gradient of a long term interaction has exponentially smaller magnitude than the gradient of a short term interaction. It does not mean that it is impossible to learn, but that it might take a very long time to learn long-term dependencies, because the signal about these dependencies will tend to be hidden by the smallest fluctuations arising from short-term dependencies. In practice, the experiments in [Bengio et al. \(1994\)](#) show that as we increase the span of the dependencies that

need to be captured, gradient-based optimization becomes increasingly difficult, with the probability of successful training of a traditional RNN via SGD rapidly reaching 0 for sequences of only length 10 or 20.

For a deeper treatment of the dynamical systems view of recurrent networks, see Doya (1993), Bengio *et al.* (1994) and Siegelmann and Sontag (1995), with a review in Pascanu *et al.* (2013a). The remaining sections of this chapter discuss various approaches that have been proposed to reduce the difficulty of learning long-term dependencies (in some cases allowing an RNN to learn dependencies across hundreds of steps), but the problem of learning long-term dependencies remains one of the main challenges in deep learning.

## 10.8 Echo State Networks

The recurrent weights mapping from  $\mathbf{h}^{(t-1)}$  to  $\mathbf{h}^{(t)}$  and the input weights mapping from  $\mathbf{x}^{(t)}$  to  $\mathbf{h}^{(t)}$  are some of the most difficult parameters to learn in a recurrent network. One proposed (Jaeger, 2003; Maass *et al.*, 2002; Jaeger and Haas, 2004; Jaeger, 2007b) approach to avoiding this difficulty is to set the recurrent weights such that the recurrent hidden units do a good job of capturing the history of past inputs, and **only learn the output weights**. This is the idea that was independently proposed for *echo state networks* or ESNs (Jaeger and Haas, 2004; Jaeger, 2007b) and *liquid state machines* (Maass *et al.*, 2002). The latter is similar, except that it uses spiking neurons (with binary outputs) instead of the continuous-valued hidden units used for ESNs. Both ESNs and liquid state machines are termed *reservoir computing* (Lukoševičius and Jaeger, 2009) to denote the fact that the hidden units form of reservoir of temporal features which may capture different aspects of the history of inputs.

One way to think about these reservoir computing recurrent networks is that they are similar to kernel machines: they map an arbitrary length sequence (the history of inputs up to time  $t$ ) into a fixed-length vector (the recurrent state  $\mathbf{h}^{(t)}$ ), on which a linear predictor (typically a linear regression) can be applied to solve the problem of interest. The training criterion may then be easily designed to be convex as a function of the output weights. For example, if the output consists of linear regression from the hidden units to the output targets, and the training criterion is mean squared error, then it is convex and may be solved reliably with simple learning algorithms (Jaeger, 2003).

The important question is therefore: how do we set the input and recurrent weights so that a rich set of histories can be represented in the recurrent neural network state? The answer proposed in the reservoir computing literature is to

view the recurrent net as a dynamical system, and set the input and recurrent weights such that the dynamical system is near the edge of stability.

The original idea was to make the eigenvalues of the Jacobian of the state-to-state transition function be close to 1. As explained in Sec. 8.2.5, an important characteristic of a recurrent network is the eigenvalue spectrum of the Jacobians  $\mathbf{J}^{(t)} = \frac{\partial s^{(t)}}{\partial s^{(t-1)}}$ . Of particular importance is the *spectral radius* of  $\mathbf{J}^{(t)}$ , defined to be the maximum of the absolute values of its eigenvalues.

To understand the effect of the spectral radius, consider the simple case of back-propagation with a Jacobian matrix  $\mathbf{J}$  that does not change with  $t$ . This case happens, for example, when the network is purely linear. Suppose that  $\mathbf{J}$  has an eigenvector  $\mathbf{v}$  with corresponding eigenvalue  $\lambda$ . Consider what happens as we propagate a gradient vector backwards through time. If we begin with a gradient vector  $\mathbf{g}$ , then after one step of back-propagation, we will have  $\mathbf{J}\mathbf{g}$ , and after  $n$  steps we will have  $\mathbf{J}^n\mathbf{g}$ . Now consider what happens if we instead back-propagate a perturbed version of  $\mathbf{g}$ . If we begin with  $\mathbf{g} + \delta\mathbf{v}$ , then after one step, we will have  $\mathbf{J}(\mathbf{g} + \delta\mathbf{v})$ . After  $n$  steps, we will have  $\mathbf{J}^n(\mathbf{g} + \delta\mathbf{v})$ . From this we can see that back-propagation starting from  $\mathbf{g}$  and back-propagation starting from  $\mathbf{g} + \delta\mathbf{v}$  diverge by  $\delta\mathbf{J}^n\mathbf{v}$  after  $n$  steps of back-propagation. If  $\mathbf{v}$  is chosen to be a unit eigenvector of  $\mathbf{J}$  with eigenvalue  $\lambda$ , then multiplication by the Jacobian simply scales the difference at each step. The two executions of back-propagation are separated by a distance of  $\delta|\lambda|^n$ . When  $\mathbf{v}$  corresponds to the largest value of  $|\lambda|$ , this perturbation achieves the widest possible separation of an initial perturbation of size  $\delta$ .

When  $|\lambda| > 1$ , the deviation size  $\delta|\lambda|^n$  grows exponentially large. When  $|\lambda| < 1$ , the deviation size becomes exponentially small.

Of course, this example assumed that the Jacobian was the same at every time step, corresponding to a recurrent network with no nonlinearity. When a nonlinearity is present, the derivative of the nonlinearity will approach zero on many time steps, and help to prevent the explosion resulting from a large spectral radius. Indeed, the most recent work on echo state networks advocates using a spectral radius much larger than unity (Yildiz *et al.*, 2012; Jaeger, 2012).

Everything we have said about back-propagation via repeated matrix multiplication applies equally to forward propagation in a network with no nonlinearity, where the state  $\mathbf{h}^{(t+1)} = \mathbf{h}^{(t)\top} \mathbf{W}$ .

When a linear map  $\mathbf{W}^\top$  always shrinks  $\mathbf{h}$  as measured by the  $L^2$  norm, then we say that the map is *contractive*. When the spectral radius is less than one, the mapping from  $\mathbf{h}^{(t)}$  to  $\mathbf{h}^{(t+1)}$  is contractive, so a small change becomes smaller after each time step. This necessarily makes the network forget information about the

past when we use a finite level of precision (such as 32 bit integers) to store the state vector.

The Jacobian matrix tells us how a small change of  $\mathbf{h}^{(t)}$  propagates one step forward, or equivalently, how the gradient on  $\mathbf{h}^{(t+1)}$  propagates one step backward, during back-propagation. Note that neither  $\mathbf{W}$  nor  $\mathbf{J}$  need to be symmetric (although they are square and real), so they can have complex-valued eigenvalues and eigenvectors, with imaginary components corresponding to potentially oscillatory behavior (if the same Jacobian was applied iteratively). Even though  $\mathbf{h}^{(t)}$  or a small variation of  $\mathbf{h}^{(t)}$  of interest in back-propagation are real-valued, they can be expressed in such a complex-valued basis. What matters is what happens to the magnitude (complex absolute value) of these possibly complex-valued basis coefficients, when we multiply the matrix by the vector. An eigenvalue with magnitude greater than one corresponds to magnification (exponential growth, if applied iteratively) or shrinking (exponential decay, if applied iteratively).

With a nonlinear map, the Jacobian is free to change at each step. The dynamics therefore become more complicated. However, it remains true that a small initial variation can turn into a large variation after several steps. One difference between the purely linear case and the nonlinear case is that the use of a squashing nonlinearity such as tanh can cause the recurrent dynamics to become bounded. Note that it is possible for back-propagation to retain unbounded dynamics even when forward propagation has bounded dynamics, for example, when a sequence of tanh units are all in the middle of their linear regime and are connected by weight matrices with spectral radius greater than 1. However, it is rare for all of the tanh units to simultaneously lie at their linear activation point.

The strategy of echo state networks is simply to fix the weights to have some spectral radius such as 3, where information is carried forward through time but does not explode due to the stabilizing effect of saturating nonlinearities like tanh.

More recently, it has been shown that the techniques used to set the weights in ESNs could be used to **initialize** the weights in a fully trainable recurrent network (with the hidden-to-hidden recurrent weights trained using back-propagation through time), helping to learn long-term dependencies ([Sutskever, 2012](#); [Sutskever et al., 2013](#)). In this setting, an initial spectral radius of 1.2 performs well, combined with the sparse initialization scheme described in Sec. 8.4.

## 10.9 Leaky Units and Other Strategies for Multiple Time Scales

One way to deal with long-term dependencies is to design a model that operates at multiple time scales, so that some parts of the model operate at fine-grained time scales and can handle small details, while other parts operate at coarse time scales and transfer information from the distant past to the present more efficiently. Various strategies for building both fine and coarse time scales are possible. These include the addition of skip connections across time, “leaky units” that integrate signals with different time constants, and the removal of some of the connections used to model fine-grained time scales.

### 10.9.1 Adding Skip Connections through Time

One way to obtain coarse time scales is to add direct connections from variables in the distant past to variables in the present. The idea of using such skip connections dates back to Lin *et al.* (1996) and follows from the idea of incorporating delays in feedforward neural networks (Lang and Hinton, 1988). In an ordinary recurrent network, a recurrent connection goes from a unit at time  $t$  to a unit at time  $t+1$ . It is possible to construct recurrent networks with longer delays (Bengio, 1991).

As we have seen in Sec. 8.2.5, gradients may vanish or explode exponentially **with respect to the number of time steps**. Lin *et al.* (1996) introduced recurrent connections with a time-delay of  $d$  to mitigate this problem. Gradients now diminish exponentially as a function of  $\frac{\tau}{d}$  rather than  $\tau$ . Since there are both delayed and single step connections, gradients may still explode exponentially in  $\tau$ . This allows the learning algorithm to capture longer dependencies although not all long-term dependencies may be represented well in this way.

### 10.9.2 Leaky Units and a Spectrum of Different Time Scales

Another way to obtain paths on which the product of derivatives is close to one is to have units with **linear** self-connections and a weight near one on these connections.

When we accumulate a running average  $\mu^{(t)}$  of some value  $v^{(t)}$  by applying the update  $\mu^{(t)} \leftarrow \alpha\mu^{(t-1)} + (1 - \alpha)v^{(t)}$  the  $\alpha$  parameter is an example of a linear self-connection from  $\mu^{(t-1)}$  to  $\mu^{(t)}$ . When  $\alpha$  is near one, the running average remembers information about the past for a long time, and when  $\alpha$  is near zero, information about the past is rapidly discarded. Hidden units with linear self-connections can behave similarly to such running averages. Such hidden units are called *leaky units*.

Skip connections through  $d$  time steps are a way of ensuring that a unit can always learn to be influenced by a value from  $d$  time steps earlier. The use of a linear self-connection with a weight near one is a different way of ensuring that the unit can access values from the past. The linear self-connection approach allows this effect to be adapted more smoothly and flexibly by adjusting the real-valued  $\alpha$  rather than by adjusting the integer-valued skip length.

These ideas were proposed by Mozer (1992) and by El Hihi and Bengio (1996). Leaky units were also found to be useful in the context of echo state networks (Jaeger *et al.*, 2007).

There are two basic strategies for setting the time constants used by leaky units. One strategy is to manually fix them to values that remain constant, for example by sampling their values from some distribution once at initialization time. Another strategy is to make the time constants free parameters and learn them. Having such leaky units at different time scales appears to help with long-term dependencies (Mozer, 1992; Pascanu *et al.*, 2013a).

### 10.9.3 Removing Connections

Another approach to handle long-term dependencies is the idea of organizing the state of the RNN at multiple time-scales (El Hihi and Bengio, 1996), with information flowing more easily through long distances at the slower time scales.

This idea differs from the skip connections through time discussed earlier because it involves actively **removing** length-one connections and replacing them with longer connections. Units modified in such a way are forced to operate on a long time scale. Skip connections through time **add** edges. Units receiving such new connections may learn to operate on a long time scale but may also choose to focus on their other short-term connections.

There are different ways in which a group of recurrent units can be forced to operate at different time scales. One option is to make the recurrent units leaky, but to have different groups of units associated with different fixed time scales. This was the proposal in Mozer (1992) and has been successfully used in Pascanu *et al.* (2013a). Another option is to have explicit and discrete updates taking place at different times, with a different frequency for different groups of units. This is the approach of El Hihi and Bengio (1996) and Koutnik *et al.* (2014). It worked well on a number of benchmark datasets.

## 10.10 The Long Short-Term Memory and Other Gated RNNs

As of this writing, the most effective sequence models used in practical applications are called *gated RNNs*. These include the *long short-term memory* and networks based on the *gated recurrent unit*.

Like leaky units, gated RNNs are based on the idea of creating paths through time that have derivatives that neither vanish nor explode. Leaky units did this with connection weights that were either manually chosen constants or were parameters. Gated RNNs generalize this to connection weights that may change at each time step.

Leaky units allow the network to **accumulate** information (such as evidence for a particular feature or category) over a long duration. However, once that information has been used, it might be useful for the neural network to **forget** the old state. For example, if a sequence is made of sub-sequences and we want a leaky unit to accumulate evidence inside each sub-subsequence, we need a mechanism to forget the old state by setting it to zero. Instead of manually deciding when to clear the state, we want the neural network to learn to decide when to do it. This is what gated RNNs do.

### 10.10.1 LSTM

The clever idea of introducing self-loops to produce paths where the gradient can flow for long durations is a core contribution of the initial *long short-term memory* (*LSTM*) model (Hochreiter and Schmidhuber, 1997). A crucial addition has been to make the weight on this self-loop conditioned on the context, rather than fixed (Gers *et al.*, 2000). By making the weight of this self-loop gated (controlled by another hidden unit), the time scale of integration can be changed dynamically. In this case, we mean that even for an LSTM with fixed parameters, the time scale of integration can change based on the input sequence, because the time constants are output by the model itself. The LSTM has been found extremely successful in many applications, such as unconstrained handwriting recognition (Graves *et al.*, 2009), speech recognition (Graves *et al.*, 2013; Graves and Jaitly, 2014), handwriting generation (Graves, 2013), machine translation (Sutskever *et al.*, 2014), image captioning (Kiros *et al.*, 2014b; Vinyals *et al.*, 2014b; Xu *et al.*, 2015) and parsing (Vinyals *et al.*, 2014a).

The LSTM block diagram is illustrated in Fig. 10.16. The corresponding forward propagation equations are given below, in the case of a shallow recurrent

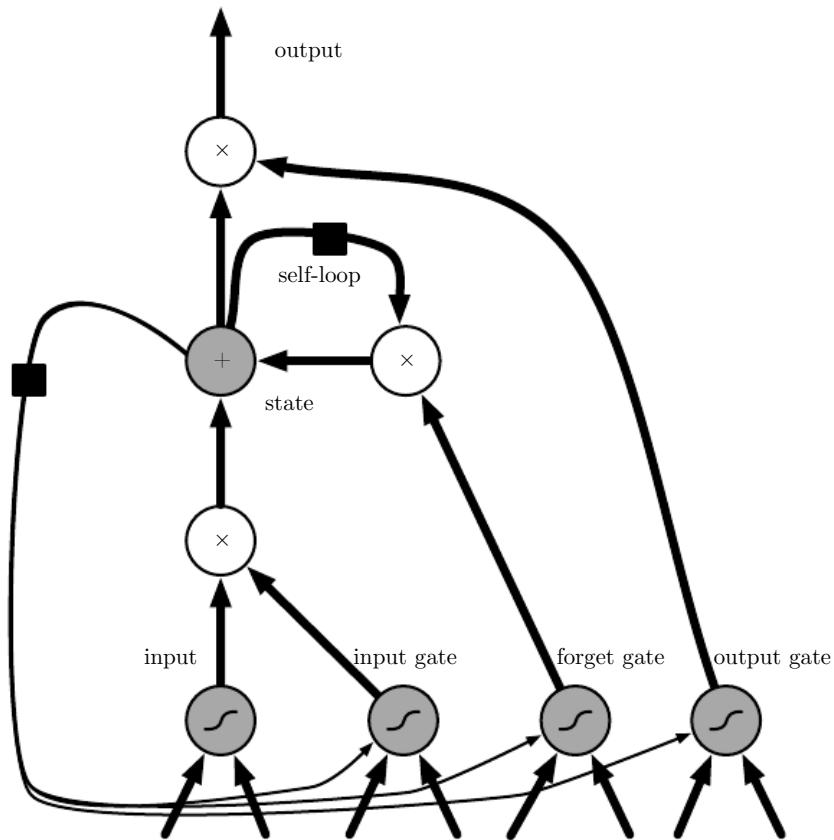


Figure 10.16: Block diagram of the LSTM recurrent network “cell.” Cells are connected recurrently to each other, replacing the usual hidden units of ordinary recurrent networks. An input feature is computed with a regular artificial neuron unit. Its value can be accumulated into the state if the sigmoidal input gate allows it. The state unit has a linear self-loop whose weight is controlled by the forget gate. The output of the cell can be shut off by the output gate. All the gating units have a sigmoid nonlinearity, while the input unit can have any squashing nonlinearity. The state unit can also be used as an extra input to the gating units. The black square indicates a delay of 1 time unit.

network architecture. Deeper architectures have also been successfully used (Graves *et al.*, 2013; Pascanu *et al.*, 2014a). Instead of a unit that simply applies an element-wise nonlinearity to the affine transformation of inputs and recurrent units, LSTM recurrent networks have “LSTM cells” that have an internal recurrence (a self-loop), in addition to the outer recurrence of the RNN. Each cell has the same inputs and outputs as an ordinary recurrent network, but has more parameters and a system of gating units that controls the flow of information. The most important component is the state unit  $s_i^{(t)}$  that has a linear self-loop similar to the leaky units described in the previous section. However, here, the self-loop weight (or the associated time constant) is controlled by a *forget gate* unit  $f_i^{(t)}$  (for time step  $t$  and cell  $i$ ), that sets this weight to a value between 0 and 1 via a sigmoid unit:

$$f_i^{(t)} = \sigma \left( b_i^f + \sum_j U_{i,j}^f x_j^{(t)} + \sum_j W_{i,j}^f h_j^{(t-1)} \right), \quad (10.33)$$

where  $\mathbf{x}^{(t)}$  is the current input vector and  $\mathbf{h}^{(t)}$  is the current hidden layer vector, containing the outputs of all the LSTM cells, and  $\mathbf{b}^f, \mathbf{U}^f, \mathbf{W}^f$  are respectively biases, input weights and recurrent weights for the forget gates. The LSTM cell internal state is thus updated as follows, but with a conditional self-loop weight  $f_i^{(t)}$ :

$$s_i^{(t)} = f_i^{(t)} s_i^{(t-1)} + g_i^{(t)} \sigma \left( b_i + \sum_j U_{i,j} x_j^{(t)} + \sum_j W_{i,j} h_j^{(t-1)} \right), \quad (10.34)$$

where  $\mathbf{b}, \mathbf{U}$  and  $\mathbf{W}$  respectively denote the biases, input weights and recurrent weights into the LSTM cell. The *external input gate* unit  $g_i^{(t)}$  is computed similarly to the forget gate (with a sigmoid unit to obtain a gating value between 0 and 1), but with its own parameters:

$$g_i^{(t)} = \sigma \left( b_i^g + \sum_j U_{i,j}^g x_j^{(t)} + \sum_j W_{i,j}^g h_j^{(t-1)} \right). \quad (10.35)$$

The output  $h_i^{(t)}$  of the LSTM cell can also be shut off, via the *output gate*  $q_i^{(t)}$ , which also uses a sigmoid unit for gating:

$$h_i^{(t)} = \tanh(s_i^{(t)}) q_i^{(t)} \quad (10.36)$$

$$q_i^{(t)} = \sigma \left( b_i^o + \sum_j U_{i,j}^o x_j^{(t)} + \sum_j W_{i,j}^o h_j^{(t-1)} \right) \quad (10.37)$$

which has parameters  $\mathbf{b}^o$ ,  $\mathbf{U}^o$ ,  $\mathbf{W}^o$  for its biases, input weights and recurrent weights, respectively. Among the variants, one can choose to use the cell state  $s_i^{(t)}$  as an extra input (with its weight) into the three gates of the  $i$ -th unit, as shown in Fig. 10.16. This would require three additional parameters.

LSTM networks have been shown to learn long-term dependencies more easily than the simple recurrent architectures, first on artificial data sets designed for testing the ability to learn long-term dependencies (Bengio *et al.*, 1994; Hochreiter and Schmidhuber, 1997; Hochreiter *et al.*, 2000), then on challenging sequence processing tasks where state-of-the-art performance was obtained (Graves, 2012; Graves *et al.*, 2013; Sutskever *et al.*, 2014). Variants and alternatives to the LSTM have been studied and used and are discussed next.

### 10.10.2 Other Gated RNNs

Which pieces of the LSTM architecture are actually necessary? What other successful architectures could be designed that allow the network to dynamically control the time scale and forgetting behavior of different units?

Some answers to these questions are given with the recent work on gated RNNs, whose units are also known as gated recurrent units or GRUs (Cho *et al.*, 2014b; Chung *et al.*, 2014, 2015a; Jozefowicz *et al.*, 2015; Chrupala *et al.*, 2015). The main difference with the LSTM is that a single gating unit simultaneously controls the forgetting factor and the decision to update the state unit. The update equations are the following:

$$h_i^{(t)} = u_i^{(t-1)} h_i^{(t-1)} + (1 - u_i^{(t-1)}) \sigma \left( b_i + \sum_j U_{i,j} x_j^{(t-1)} + \sum_j W_{i,j} r_j^{(t-1)} h_j^{(t-1)} \right), \quad (10.38)$$

where  $\mathbf{u}$  stands for “update” gate and  $\mathbf{r}$  for “reset” gate. Their value is defined as usual:

$$u_i^{(t)} = \sigma \left( b_i^u + \sum_j U_{i,j}^u x_j^{(t)} + \sum_j W_{i,j}^u h_j^{(t)} \right) \quad (10.39)$$

and

$$r_i^{(t)} = \sigma \left( b_i^r + \sum_j U_{i,j}^r x_j^{(t)} + \sum_j W_{i,j}^r h_j^{(t)} \right). \quad (10.40)$$

The reset and updates gates can individually “ignore” parts of the state vector. The update gates act like conditional leaky integrators that can linearly gate any

dimension, thus choosing to copy it (at one extreme of the sigmoid) or completely ignore it (at the other extreme) by replacing it by the new “target state” value (towards which the leaky integrator wants to converge). The reset gates control which parts of the state get used to compute the next target state, introducing an additional nonlinear effect in the relationship between past state and future state.

Many more variants around this theme can be designed. For example the reset gate (or forget gate) output could be shared across multiple hidden units. Alternately, the product of a global gate (covering a whole group of units, such as an entire layer) and a local gate (per unit) could be used to combine global control and local control. However, several investigations over architectural variations of the LSTM and GRU found no variant that would clearly beat both of these across a wide range of tasks (Greff *et al.*, 2015; Jozefowicz *et al.*, 2015). Greff *et al.* (2015) found that a crucial ingredient is the forget gate, while Jozefowicz *et al.* (2015) found that adding a bias of 1 to the LSTM forget gate, a practice advocated by Gers *et al.* (2000), makes the LSTM as strong as the best of the explored architectural variants.

## 10.11 Optimization for Long-Term Dependencies

Sec. 8.2.5 and Sec. 10.7 have described the vanishing and exploding gradient problems that occur when optimizing RNNs over many time steps.

An interesting idea proposed by Martens and Sutskever (2011) is that second derivatives may vanish at the same time that first derivatives vanish. Second-order optimization algorithms may roughly be understood as dividing the first derivative by the second derivative (in higher dimension, multiplying the gradient by the inverse Hessian). If the second derivative shrinks at a similar rate to the first derivative, then the ratio of first and second derivatives may remain relatively constant. Unfortunately, second-order methods have many drawbacks, including high computational cost, the need for a large minibatch, and a tendency to be attracted to saddle points. Martens and Sutskever (2011) found promising results using second-order methods. Later, Sutskever *et al.* (2013) found that simpler methods such as Nesterov momentum with careful initialization could achieve similar results. See Sutskever (2012) for more detail. Both of these approaches have largely been replaced by simply using SGD (even without momentum) applied to LSTMs. This is part of a continuing theme in machine learning that it is often much easier to design a model that is easy to optimize than it is to design a more powerful optimization algorithm.

### 10.11.1 Clipping Gradients

As discussed in Sec. 8.2.4, strongly nonlinear functions such as those computed by a recurrent net over many time steps tend to have derivatives that can be either very large or very small in magnitude. This is illustrated in Fig. 8.3 and Fig. 10.17, in which we see that the objective function (as a function of the parameters) has a “landscape” in which one finds “cliffs”: wide and rather flat regions separated by tiny regions where the objective function changes quickly, forming a kind of cliff.

The difficulty that arises is that when the parameter gradient is very large, a gradient descent parameter update could throw the parameters very far, into a region where the objective function is larger, undoing much of the work that had been done to reach the current solution. The gradient tells us the direction that corresponds to the steepest descent within an infinitesimal region surrounding the current parameters. Outside of this infinitesimal region, the cost function may begin to curve back upwards. The update must be chosen to be small enough to avoid traversing too much upward curvature. We typically use learning rates that decay slowly enough that consecutive steps have approximately the same learning rate. A step size that is appropriate for a relatively linear part of the landscape is often inappropriate and causes uphill motion if we enter a more curved part of the landscape on the next step.

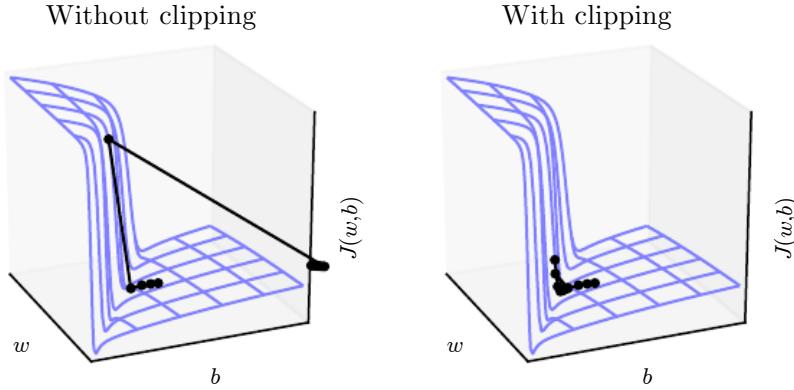


Figure 10.17: Example of the effect of gradient clipping in a recurrent network with two parameters  $\mathbf{w}$  and  $\mathbf{b}$ . Gradient clipping can make gradient descent perform more reasonably in the vicinity of extremely steep cliffs. These steep cliffs commonly occur in recurrent networks near where a recurrent network behaves approximately linearly. The cliff is exponentially steep in the number of time steps because the weight matrix is multiplied by itself once for each time step. (*Left*) Gradient descent without gradient clipping overshoots the bottom of this small ravine, then receives a very large gradient from the cliff face. The large gradient catastrophically propels the parameters outside the axes of the plot. (*Right*) Gradient descent with gradient clipping has a more moderate reaction to the cliff. While it does ascend the cliff face, the step size is restricted so that it cannot be propelled away from steep region near the solution. Figure adapted with permission from [Pascanu et al. \(2013a\)](#).

A simple type of solution has been in use by practitioners for many years: *clipping the gradient*. There are different instances of this idea ([Mikolov, 2012](#); [Pascanu et al., 2013a](#)). One option is to clip the parameter gradient from a minibatch *element-wise* ([Mikolov, 2012](#)) just before the parameter update. Another is to *clip the norm*  $\|\mathbf{g}\|$  of the gradient  $\mathbf{g}$  ([Pascanu et al., 2013a](#)) just before the parameter update:

$$\text{if } \|\mathbf{g}\| > v \quad (10.41)$$

$$\mathbf{g} \leftarrow \frac{\mathbf{g}v}{\|\mathbf{g}\|} \quad (10.42)$$

where  $v$  is the norm threshold and  $\mathbf{g}$  is used to update parameters. Because the gradient of all the parameters (including different groups of parameters, such as weights and biases) is renormalized jointly with a single scaling factor, the latter method has the advantage that it guarantees that each step is still in the gradient direction, but experiments suggest that both forms work similarly. Although

the parameter update has the same direction as the true gradient, with gradient norm clipping, the parameter update vector norm is now bounded. This bounded gradient avoids performing a detrimental step when the gradient explodes. In fact, even simply taking a *random step* when the gradient magnitude is above a threshold tends to work almost as well. If the explosion is so severe that the gradient is numerically `Inf` or `Nan` (considered infinite or not-a-number), then a random step of size  $v$  can be taken and will typically move away from the numerically unstable configuration. Clipping the gradient norm per-minibatch will not change the direction of the gradient for an individual minibatch. However, taking the average of the norm-clipped gradient from many minibatches is not equivalent to clipping the norm of the true gradient (the gradient formed from using all examples). Examples that have large gradient norm, as well as examples that appear in the same minibatch as such examples, will have their contribution to the final direction diminished. This stands in contrast to traditional minibatch gradient descent, where the true gradient direction is equal to the average over all minibatch gradients. Put another way, traditional stochastic gradient descent uses an unbiased estimate of the gradient, while gradient descent with norm clipping introduces a heuristic bias that we know empirically to be useful. With element-wise clipping, the direction of the update is not aligned with the true gradient or the minibatch gradient, but it is still a descent direction. It has also been proposed (Graves, 2013) to clip the back-propagated gradient (with respect to hidden units) but no comparison has been published between these variants; we conjecture that all these methods behave similarly.

### 10.11.2 Regularizing to Encourage Information Flow

Gradient clipping helps to deal with exploding gradients, but it does not help with vanishing gradients. To address vanishing gradients and better capture long-term dependencies, we discussed the idea of creating paths in the computational graph of the unfolded recurrent architecture along which the product of gradients associated with arcs is near 1. One approach to achieve this is with LSTMs and other self-loops and gating mechanisms, described above in Sec. 10.10. Another idea is to regularize or constrain the parameters so as to encourage “information flow.” In particular, we would like the gradient vector  $\nabla_{\mathbf{h}^{(t)}} L$  being back-propagated to maintain its magnitude, even if the loss function only penalizes the output at the end of the sequence. Formally, we want

$$(\nabla_{\mathbf{h}^{(t)}} L) \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} \tag{10.43}$$

to be as large as

$$\nabla_{\mathbf{h}^{(t)}} L. \quad (10.44)$$

With this objective, [Pascanu et al. \(2013a\)](#) propose the following regularizer:

$$\Omega = \sum_t \left( \frac{\left| \left| (\nabla_{\mathbf{h}^{(t)}} L) \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} \right| \right|}{\|\nabla_{\mathbf{h}^{(t)}} L\|} - 1 \right)^2. \quad (10.45)$$

Computing the gradient of this regularizer may appear difficult, but [Pascanu et al. \(2013a\)](#) propose an approximation in which we consider the back-propagated vectors  $\nabla_{\mathbf{h}^{(t)}} L$  as if they were constants (for the purpose of this regularizer, so that there is no need to back-propagate through them). The experiments with this regularizer suggest that, if combined with the norm clipping heuristic (which handles gradient explosion), the regularizer can considerably increase the span of the dependencies that an RNN can learn. Because it keeps the RNN dynamics on the edge of explosive gradients, the gradient clipping is particularly important. Without gradient clipping, gradient explosion prevents learning from succeeding.

A key weakness of this approach is that it is not as effective as the LSTM for tasks where data is abundant, such as language modeling.

## 10.12 Explicit Memory

Intelligence requires knowledge and acquiring knowledge can be done via learning, which has motivated the development of large-scale deep architectures. However, there are different kinds of knowledge. Some knowledge can be implicit, subconscious, and difficult to verbalize—such as how to walk, or how a dog looks different from a cat. Other knowledge can be explicit, declarative, and relatively straightforward to put into words—every day commonsense knowledge, like “a cat is a kind of animal,” or very specific facts that you need to know to accomplish your current goals, like “the meeting with the sales team is at 3:00 PM in room 141.”

Neural networks excel at storing implicit knowledge. However, they struggle to memorize facts. Stochastic gradient descent requires many presentations of the same input before it can be stored in a neural network parameters, and even then, that input will not be stored especially precisely. [Graves et al. \(2014b\)](#) hypothesized that this is because neural networks lack the equivalent of the *working memory* system that allows human beings to explicitly hold and manipulate pieces of information that are relevant to achieving some goal. Such explicit memory

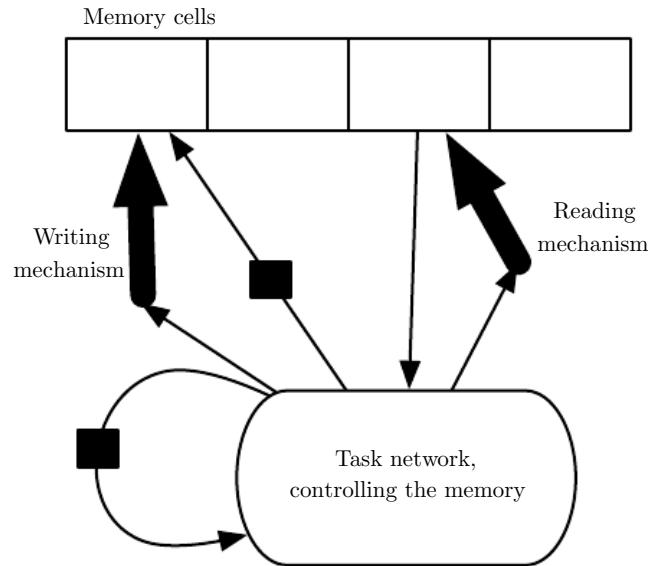


Figure 10.18: A schematic of an example of a network with an explicit memory, capturing some of the key design elements of the neural Turing machine. In this diagram we distinguish the “representation” part of the model (the “task network,” here a recurrent net in the bottom) from the “memory” part of the model (the set of cells), which can store facts. The task network learns to “control” the memory, deciding where to read from and where to write to within the memory (through the reading and writing mechanisms, indicated by bold arrows pointing at the reading and writing addresses).

components would allow our systems not only to rapidly and “intentionally” store and retrieve specific facts but also to sequentially reason with them. The need for neural networks that can process information in a sequence of steps, changing the way the input is fed into the network at each step, has long been recognized as important for the ability to reason rather than to make automatic, intuitive responses to the input (Hinton, 1990).

To resolve this difficulty, Weston *et al.* (2014) introduced *memory networks* that include a set of memory cells that can be accessed via an addressing mechanism. Memory networks originally required a supervision signal instructing them how to use their memory cells. Graves *et al.* (2014b) introduced the *neural Turing machine*, which is able to learn to read from and write arbitrary content to memory cells without explicit supervision about which actions to undertake, and allowed end-to-end training without this supervision signal, via the use of a content-based soft attention mechanism (see Bahdanau *et al.* (2015) and Sec. 12.4.5.1). This soft addressing mechanism has become standard with other related architectures emulating algorithmic mechanisms in a way that still allows gradient-based optimization (Sukhbaatar *et al.*, 2015; Joulin and Mikolov, 2015; Kumar *et al.*, 2015; Vinyals *et al.*, 2015a; Grefenstette *et al.*, 2015).

Each memory cell can be thought of as an extension of the memory cells in LSTMs and GRUs. The difference is that the network outputs an internal state that chooses which cell to read from or write to, just as memory accesses in a digital computer read from or write to a specific address.

It is difficult to optimize functions that produce exact, integer addresses. To alleviate this problem, NTMs actually read to or write from many memory cells simultaneously. To read, they take a weighted average of many cells. To write, they modify multiple cells by different amounts. The coefficients for these operations are chosen to be focused on a small number of cells, for example, by producing them via a softmax function. Using these weights with non-zero derivatives allows the functions controlling access to the memory to be optimized using gradient descent. The gradient on these coefficients indicates whether each of them should be increased or decreased, but the gradient will typically be large only for those memory addresses receiving a large coefficient.

These memory cells are typically augmented to contain a vector, rather than the single scalar stored by an LSTM or GRU memory cell. There are two reasons to increase the size of the memory cell. One reason is that we have increased the cost of accessing a memory cell. We pay the computational cost of producing a coefficient for many cells, but we expect these coefficients to cluster around a small number of cells. By reading a vector value, rather than a scalar value, we can

offset some of this cost. Another reason to use vector-valued memory cells is that they allow for *content-based addressing*, where the weight used to read to or write from a cell is a function of that cell. Vector-valued cells allow us to retrieve a complete vector-valued memory if we are able to produce a pattern that matches some but not all of its elements. This is analogous to the way that people can recall the lyrics of a song based on a few words. We can think of a content-based read instruction as saying, “Retrieve the lyrics of the song that has the chorus ‘We all live in a yellow submarine.’” Content-based addressing is more useful when we make the objects to be retrieved large—if every letter of the song was stored in a separate memory cell, we would not be able to find them this way. By comparison, *location-based addressing* is not allowed to refer to the content of the memory. We can think of a location-based read instruction as saying “Retrieve the lyrics of the song in slot 347.” Location-based addressing can often be a perfectly sensible mechanism even when the memory cells are small.

If the content of a memory cell is copied (not forgotten) at most time steps, then the information it contains can be propagated forward in time and the gradients propagated backward in time without either vanishing or exploding.

The explicit memory approach is illustrated in Fig. 10.18, where we see that a “task neural network” is coupled with a memory. Although that task neural network could be feedforward or recurrent, the overall system is a recurrent network. The task network can choose to read from or write to specific memory addresses. Explicit memory seems to allow models to learn tasks that ordinary RNNs or LSTM RNNs cannot learn. One reason for this advantage may be because information and gradients can be propagated (forward in time or backwards in time, respectively) for very long durations.

As an alternative to back-propagation through weighted averages of memory cells, we can interpret the memory addressing coefficients as probabilities and stochastically read just one cell (Zaremba and Sutskever, 2015). Optimizing models that make discrete decisions requires specialized optimization algorithms, described in Sec. 20.9.1. So far, training these stochastic architectures that make discrete decisions remains harder than training deterministic algorithms that make soft decisions.

Whether it is soft (allowing back-propagation) or stochastic and hard, the mechanism for choosing an address is in its form identical to the *attention mechanism* which had been previously introduced in the context of machine translation (Bahdanau *et al.*, 2015) and discussed in Sec. 12.4.5.1. The idea of attention mechanisms for neural networks was introduced even earlier, in the context of handwriting generation (Graves, 2013), with an attention mechanism that was constrained to

move only forward in time through the sequence. In the case of machine translation and memory networks, at each step, the focus of attention can move to a completely different place, compared to the previous step.

Recurrent neural networks provide a way to extend deep learning to sequential data. They are the last major tool in our deep learning toolbox. Our discussion now moves to how to choose and use these tools and how to apply them to real-world tasks.

# Chapter 11

## Practical methodology

Successfully applying deep learning techniques requires more than just a good knowledge of what algorithms exist and the principles that explain how they work. A good machine learning practitioner also needs to know how to choose an algorithm for a particular application and how to monitor and respond to feedback obtained from experiments in order to improve a machine learning system. During day to day development of machine learning systems, practitioners need to decide whether to gather more data, increase or decrease model capacity, add or remove regularizing features, improve the optimization of a model, improve approximate inference in a model, or debug the software implementation of the model. All of these operations are at the very least time-consuming to try out, so it is important to be able to determine the right course of action rather than blindly guessing.

Most of this book is about different machine learning models, training algorithms, and objective functions. This may give the impression that the most important ingredient to being a machine learning expert is knowing a wide variety of machine learning techniques and being good at different kinds of math. In practice, one can usually do much better with a correct application of a commonplace algorithm than by sloppily applying an obscure algorithm. Correct application of an algorithm depends on mastering some fairly simple methodology. Many of the recommendations in this chapter are adapted from Ng (2015).

We recommend the following practical design process:

- Determine your goals—what error metric to use, and your target value for this error metric. These goals and error metrics should be driven by the problem that the application is intended to solve.
- Establish a working end-to-end pipeline as soon as possible, including the

estimation of the appropriate performance metrics.

- Instrument the system well to determine bottlenecks in performance. Diagnose which components are performing worse than expected and whether it is due to overfitting, underfitting, or a defect in the data or software.
- Repeatedly make incremental changes such as gathering new data, adjusting hyperparameters, or changing algorithms, based on specific findings from your instrumentation.

As a running example, we will use Street View address number transcription system ([Goodfellow et al., 2014d](#)). The purpose of this application is to add buildings to Google Maps. Street View cars photograph the buildings and record the GPS coordinates associated with each photograph. A convolutional network recognizes the address number in each photograph, allowing the Google Maps database to add that address in the correct location. The story of how this commercial application was developed gives an example of how to follow the design methodology we advocate.

We now describe each of the steps in this process.

## 11.1 Performance Metrics

Determining your goals, in terms of which error metric to use, is a necessary first step because your error metric will guide all of your future actions. You should also have an idea of what level of performance you desire.

Keep in mind that for most applications, it is impossible to achieve absolute zero error. The Bayes error defines the minimum error rate that you can hope to achieve, even if you have infinite training data and can recover the true probability distribution. This is because your input features may not contain complete information about the output variable, or because the system might be intrinsically stochastic. You will also be limited by having a finite amount of training data.

The amount of training data can be limited for a variety of reasons. When your goal is to build the best possible real-world product or service, you can typically collect more data but must determine the value of reducing error further and weigh this against the cost of collecting more data. Data collection can require time, money, or human suffering (for example, if your data collection process involves performing invasive medical tests). When your goal is to answer a scientific question about which algorithm performs better on a fixed benchmark, the benchmark

specification usually determines the training set and you are not allowed to collect more data.

How can one determine a reasonable level of performance to expect? Typically, in the academic setting, we have some estimate of the error rate that is attainable based on previously published benchmark results. In the real-word setting, we have some idea of the error rate that is necessary for an application to be safe, cost-effective, or appealing to consumers. Once you have determined your realistic desired error rate, your design decisions will be guided by reaching this error rate.

Another important consideration besides the target value of the performance metric is the choice of which metric to use. Several different performance metrics may be used to measure the effectiveness of a complete application that includes machine learning components. These performance metrics are usually different from the cost function used to train the model. As described in Sec. 5.1.2, it is common to measure the accuracy, or equivalently, the error rate, of a system.

However, many applications require more advanced metrics.

Sometimes it is much more costly to make one kind of a mistake than another. For example, an e-mail spam detection system can make two kinds of mistakes: incorrectly classifying a legitimate message as spam, and incorrectly allowing a spam message to appear in the inbox. It is much worse to block a legitimate message than to allow a questionable message to pass through. Rather than measuring the error rate of a spam classifier, we may wish to measure some form of total cost, where the cost of blocking legitimate messages is higher than the cost of allowing spam messages.

Sometimes we wish to train a binary classifier that is intended to detect some rare event. For example, we might design a medical test for a rare disease. Suppose that only one in every million people has this disease. We can easily achieve 99.9999% accuracy on the detection task, by simply hard-coding the classifier to always report that the disease is absent. Clearly, accuracy is a poor way to characterize the performance of such a system. One way to solve this problem is to instead measure *precision* and *recall*. Precision is the fraction of detections reported by the model that were correct, while recall is the fraction of true events that were detected. A detector that says no one has the disease would achieve perfect precision, but zero recall. A detector that says everyone has the disease would achieve perfect recall, but precision equal to the percentage of people who have the disease (0.0001% in our example of a disease that only one people in a million have). When using precision and recall, it is common to plot a *PR curve*, with precision on the *y*-axis and recall on the *x*-axis. The classifier generates a score that is higher if the event to be detected occurred. For example, a feedforward

network designed to detect a disease outputs  $\hat{y} = P(y = 1 \mid \mathbf{x})$ , estimating the probability that a person whose medical results are described by features  $\mathbf{x}$  has the disease. We choose to report a detection whenever this score exceeds some threshold. By varying the threshold, we can trade precision for recall. In many cases, we wish to summarize the performance of the classifier with a single number rather than a curve. To do so, we can convert precision  $p$  and recall  $r$  into an *F-score* given by

$$F = \frac{2pr}{p + r}. \quad (11.1)$$

Another option is to report the total area lying beneath the PR curve.

In some applications, it is possible for the machine learning system to refuse to make a decision. This is useful when the machine learning algorithm can estimate how confident it should be about a decision, especially if a wrong decision can be harmful and if a human operator is able to occasionally take over. The Street View transcription system provides an example of this situation. The task is to transcribe the address number from a photograph in order to associate the location where the photo was taken with the correct address in a map. Because the value of the map degrades considerably if the map is inaccurate, it is important to add an address only if the transcription is correct. If the machine learning system thinks that it is less likely than a human being to obtain the correct transcription, then the best course of action is to allow a human to transcribe the photo instead. Of course, the machine learning system is only useful if it is able to dramatically reduce the amount of photos that the human operators must process. A natural performance metric to use in this situation is *coverage*. Coverage is the fraction of examples for which the machine learning system is able to produce a response. It is possible to trade coverage for accuracy. One can always obtain 100% accuracy by refusing to process any example, but this reduces the coverage to 0%. For the Street View task, the goal for the project was to reach human-level transcription accuracy while maintaining 95% coverage. Human-level performance on this task is 98% accuracy.

Many other metrics are possible. We can for example, measure click-through rates, collect user satisfaction surveys, and so on. Many specialized application areas have application-specific criteria as well.

What is important is to determine which performance metric to improve ahead of time, then concentrate on improving this metric. Without clearly defined goals, it can be difficult to tell whether changes to a machine learning system make progress or not.

## 11.2 Default Baseline Models

After choosing performance metrics and goals, the next step in any practical application is to establish a reasonable end-to-end system as soon as possible. In this section, we provide recommendations for which algorithms to use as the first baseline approach in various situations. Keep in mind that deep learning research progresses quickly, so better default algorithms are likely to become available soon after this writing.

Depending on the complexity of your problem, you may even want to begin without using deep learning. If your problem has a chance of being solved by just choosing a few linear weights correctly, you may want to begin with a simple statistical model like logistic regression.

If you know that your problem falls into an “AI-complete” category like object recognition, speech recognition, machine translation, and so on, then you are likely to do well by beginning with an appropriate deep learning model.

First, choose the general category of model based on the structure of your data. If you want to perform supervised learning with fixed-size vectors as input, use a feedforward network with fully connected layers. If the input has known topological structure (for example, if the input is an image), use a convolutional network. In these cases, you should begin by using some kind of piecewise linear unit (ReLUs or their generalizations like Leaky ReLUs, PreLus and maxout). If your input or output is a sequence, use a gated recurrent net (LSTM or GRU).

A reasonable choice of optimization algorithm is SGD with momentum with a decaying learning rate (popular decay schemes that perform better or worse on different problems include decaying linearly until reaching a fixed minimum learning rate, decaying exponentially, or decreasing the learning rate by a factor of 2-10 each time validation error plateaus). Another very reasonable alternative is Adam. Batch normalization can have a dramatic effect on optimization performance, especially for convolutional networks and networks with sigmoidal nonlinearities. While it is reasonable to omit batch normalization from the very first baseline, it should be introduced quickly if optimization appears to be problematic.

Unless your training set contains tens of millions of examples or more, you should include some mild forms of regularization from the start. Early stopping should be used almost universally. Dropout is an excellent regularizer that is easy to implement and compatible with many models and training algorithms. Batch normalization also sometimes reduces generalization error and allows dropout to be omitted, due to the noise in the estimate of the statistics used to normalize each variable.

If your task is similar to another task that has been studied extensively, you will probably do well by first copying the model and algorithm that is already known to perform best on the previously studied task. You may even want to copy a trained model from that task. For example, it is common to use the features from a convolutional network trained on ImageNet to solve other computer vision tasks (Girshick *et al.*, 2015).

A common question is whether to begin by using unsupervised learning, described further in Part III. This is somewhat domain specific. Some domains, such as natural language processing, are known to benefit tremendously from unsupervised learning techniques such as learning unsupervised word embeddings. In other domains, such as computer vision, current unsupervised learning techniques do not bring a benefit, except in the semi-supervised setting, when the number of labeled examples is very small (Kingma *et al.*, 2014; Rasmus *et al.*, 2015). If your application is in a context where unsupervised learning is known to be important, then include it in your first end-to-end baseline. Otherwise, only use unsupervised learning in your first attempt if the task you want to solve is unsupervised. You can always try adding unsupervised learning later if you observe that your initial baseline overfits.

### 11.3 Determining Whether to Gather More Data

After the first end-to-end system is established, it is time to measure the performance of the algorithm and determine how to improve it. Many machine learning novices are tempted to make improvements by trying out many different algorithms. However, it is often much better to gather more data than to improve the learning algorithm.

How does one decide whether to gather more data? First, determine whether the performance on the training set is acceptable. If performance on the training set is poor, the learning algorithm is not using the training data that is already available, so there is no reason to gather more data. Instead, try increasing the size of the model by adding more layers or adding more hidden units to each layer. Also, try improving the learning algorithm, for example by tuning the learning rate hyperparameter. If large models and carefully tuned optimization algorithms do not work well, then the problem might be the **quality** of the training data. The data may be too noisy or may not include the right inputs needed to predict the desired outputs. This suggests starting over, collecting cleaner data or collecting a richer set of features.

If the performance on the training set is acceptable, then measure the per-

formance on a test set. If the performance on the test set is also acceptable, then there is nothing left to be done. If test set performance is much worse than training set performance, then gathering more data is one of the most effective solutions. The key considerations are the cost and feasibility of gathering more data, the cost and feasibility of reducing the test error by other means, and the amount of data that is expected to be necessary to improve test set performance significantly. At large internet companies with millions or billions of users, it is feasible to gather large datasets, and the expense of doing so can be considerably less than the other alternatives, so the answer is almost always to gather more training data. For example, the development of large labeled datasets was one of the most important factors in solving object recognition. In other contexts, such as medical applications, it may be costly or infeasible to gather more data. A simple alternative to gathering more data is to reduce the size of the model or improve regularization, by adjusting hyperparameters such as weight decay coefficients, or by adding regularization strategies such as dropout. If you find that the gap between train and test performance is still unacceptable even after tuning the regularization hyperparameters, then gathering more data is advisable.

When deciding whether to gather more data, it is also necessary to decide how much to gather. It is helpful to plot curves showing the relationship between training set size and generalization error, like in Fig. 5.4. By extrapolating such curves, one can predict how much additional training data would be needed to achieve a certain level of performance. Usually, adding a small fraction of the total number of examples will not have a noticeable impact on generalization error. It is therefore recommended to experiment with training set sizes on a logarithmic scale, for example doubling the number of examples between consecutive experiments.

If gathering much more data is not feasible, the only other way to improve generalization error is to improve the learning algorithm itself. This becomes the domain of research and not the domain of advice for applied practitioners.

## 11.4 Selecting Hyperparameters

Most deep learning algorithms come with many hyperparameters that control many aspects of the algorithm’s behavior. Some of these hyperparameters affect the time and memory cost of running the algorithm. Some of these hyperparameters affect the quality of the model recovered by the training process and its ability to infer correct results when deployed on new inputs.

There are two basic approaches to choosing these hyperparameters: choosing them manually and choosing them automatically. Choosing the hyperparameters

manually requires understanding what the hyperparameters do and how machine learning models achieve good generalization. Automatic hyperparameter selection algorithms greatly reduce the need to understand these ideas, but they are often much more computationally costly.

### 11.4.1 Manual Hyperparameter Tuning

To set hyperparameters manually, one must understand the relationship between hyperparameters, training error, generalization error and computational resources (memory and runtime). This means establishing a solid foundation on the fundamental ideas concerning the effective capacity of a learning algorithm from Chapter 5.

The goal of manual hyperparameter search is usually to find the lowest generalization error subject to some runtime and memory budget. We do not discuss how to determine the runtime and memory impact of various hyperparameters here because this is highly platform-dependent.

The primary goal of manual hyperparameter search is to adjust the effective capacity of the model to match the complexity of the task. Effective capacity is constrained by three factors: the representational capacity of the model, the ability of the learning algorithm to successfully minimize the cost function used to train the model, and the degree to which the cost function and training procedure regularize the model. A model with more layers and more hidden units per layer has higher representational capacity—it is capable of representing more complicated functions. It can not necessarily actually learn all of these functions though, if the training algorithm cannot discover that certain functions do a good job of minimizing the training cost, or if regularization terms such as weight decay forbid some of these functions.

The generalization error typically follows a U-shaped curve when plotted as a function of one of the hyperparameters, as in Fig. 5.3. At one extreme, the hyperparameter value corresponds to low capacity, and generalization error is high because training error is high. This is the underfitting regime. At the other extreme, the hyperparameter value corresponds to high capacity, and the generalization error is high because the gap between training and test error is high. Somewhere in the middle lies the optimal model capacity, which achieves the lowest possible generalization error, by adding a medium generalization gap to a medium amount of training error.

For some hyperparameters, overfitting occurs when the value of the hyperparameter is large. The number of hidden units in a layer is one such example,

because increasing the number of hidden units increases the capacity of the model. For some hyperparameters, overfitting occurs when the value of the hyperparameter is small. For example, the smallest allowable weight decay coefficient of zero corresponds to the greatest effective capacity of the learning algorithm.

Not every hyperparameter will be able to explore the entire U-shaped curve. Many hyperparameters are discrete, such as the number of units in a layer or the number of linear pieces in a maxout unit, so it is only possible to visit a few points along the curve. Some hyperparameters are binary. Usually these hyperparameters are switches that specify whether or not to use some optional component of the learning algorithm, such as a preprocessing step that normalizes the input features by subtracting their mean and dividing by their standard deviation. These hyperparameters can only explore two points on the curve. Other hyperparameters have some minimum or maximum value that prevents them from exploring some part of the curve. For example, the minimum weight decay coefficient is zero. This means that if the model is underfitting when weight decay is zero, we can not enter the overfitting region by modifying the weight decay coefficient. In other words, some hyperparameters can only subtract capacity.

The learning rate is perhaps the most important hyperparameter. If you have time to tune only one hyperparameter, tune the learning rate. It controls the effective capacity of the model in a more complicated way than other hyperparameters—the effective capacity of the model is highest when the learning rate is **correct** for the optimization problem, not when the learning rate is especially large or especially small. The learning rate has a U-shaped curve for *training* error, illustrated in Fig. 11.1. When the learning rate is too large, gradient descent can inadvertently increase rather than decrease the training error. In the idealized quadratic case, this occurs if the learning rate is at least twice as large as its optimal value (LeCun *et al.*, 1998a). When the learning rate is too small, training is not only slower, but may become permanently stuck with a high training error. This effect is poorly understood (it would not happen for a convex loss function).

Tuning the parameters other than the learning rate requires monitoring both training and test error to diagnose whether your model is overfitting or underfitting, then adjusting its capacity appropriately.

If your error on the training set is higher than your target error rate, you have no choice but to increase capacity. If you are not using regularization and you are confident that your optimization algorithm is performing correctly, then you must add more layers to your network or add more hidden units. Unfortunately, this increases the computational costs associated with the model.

If your error on the test set is higher than than your target error rate, you can

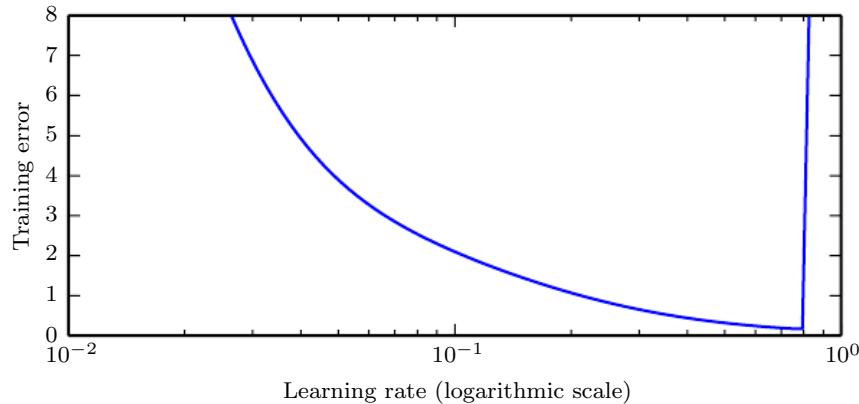


Figure 11.1: Typical relationship between the learning rate and the training error. Notice the sharp rise in error when the learning is above an optimal value. This is for a fixed training time, as a smaller learning rate may sometimes only slow down training by a factor proportional to the learning rate reduction. Generalization error can follow this curve or be complicated by regularization effects arising out of having a too large or too small learning rates, since poor optimization can, to some degree, reduce or prevent overfitting, and even points with equivalent training error can have different generalization error.

now take two kinds of actions. The test error is the sum of the training error and the gap between training and test error. The optimal test error is found by trading off these quantities. Neural networks typically perform best when the training error is very low (and thus, when capacity is high) and the test error is primarily driven by the gap between train and test error. Your goal is to reduce this gap without increasing training error faster than the gap decreases. To reduce the gap, change regularization hyperparameters to reduce effective model capacity, such as by adding dropout or weight decay. Usually the best performance comes from a large model that is regularized well, for example by using dropout.

Most hyperparameters can be set by reasoning about whether they increase or decrease model capacity. Some examples are included in Table 11.1.

While manually tuning hyperparameters, do not lose sight of your end goal: good performance on the test set. Adding regularization is only one way to achieve this goal. As long as you have low training error, you can always reduce generalization error by collecting more training data. The brute force way to practically guarantee success is to continually increase model capacity and training set size until the task is solved. This approach does of course increase the computational cost of training and inference, so it is only feasible given appropriate resources. In

Hyperparameter	Increases capacity when...	Reason	Caveats
Number of hidden units	increased	Increasing the number of hidden units increases the representational capacity of the model.	Increasing the number of hidden units increases both the time and memory cost of essentially every operation on the model.
Learning rate	tuned optimally	An improper learning rate, whether too high or too low, results in a model with low effective capacity due to optimization failure	
Convolution kernel width	increased	Increasing the kernel width increases the number of parameters in the model	A wider kernel results in a narrower output dimension, reducing model capacity unless you use implicit zero padding to reduce this effect. Wider kernels require more memory for parameter storage and increase runtime, but a narrower output reduces memory cost.
Implicit zero padding	increased	Adding implicit zeros before convolution keeps the representation size large	Increased time and memory cost of most operations.
Weight decay coefficient	decreased	Decreasing the weight decay coefficient frees the model parameters to become larger	
Dropout rate	decreased	Dropping units less often gives the units more opportunities to “conspire” with each other to fit the training set	

Table 11.1: The effect of various hyperparameters on model capacity.

principle, this approach could fail due to optimization difficulties, but for many problems optimization does not seem to be a significant barrier, provided that the model is chosen appropriately.

### 11.4.2 Automatic Hyperparameter Optimization Algorithms

The ideal learning algorithm just takes a dataset and outputs a function, without requiring hand-tuning of hyperparameters. The popularity of several learning algorithms such as logistic regression and SVMs stems in part from their ability to perform well with only one or two tuned hyperparameters. Neural networks can sometimes perform well with only a small number of tuned hyperparameters, but often benefit significantly from tuning of forty or more hyperparameters. Manual hyperparameter tuning can work very well when the user has a good starting point, such as one determined by others having worked on the same type of application and architecture, or when the user has months or years of experience in exploring hyperparameter values for neural networks applied to similar tasks. However, for many applications, these starting points are not available. In these cases, automated algorithms can find useful values of the hyperparameters.

If we think about the way in which the user of a learning algorithm searches for good values of the hyperparameters, we realize that an optimization is taking place: we are trying to find a value of the hyperparameters that optimizes an objective function, such as validation error, sometimes under constraints (such as a budget for training time, memory or recognition time). It is therefore possible, in principle, to develop *hyperparameter optimization* algorithms that wrap a learning algorithm and choose its hyperparameters, thus hiding the hyperparameters of the learning algorithm from the user. Unfortunately, hyperparameter optimization algorithms often have their own hyperparameters, such as the range of values that should be explored for each of the learning algorithm's hyperparameters. However, these secondary hyperparameters are usually easier to choose, in the sense that acceptable performance may be achieved on a wide range of tasks using the same secondary hyperparameters for all tasks.

### 11.4.3 Grid Search

When there are three or fewer hyperparameters, the common practice is to perform *grid search*. For each hyperparameter, the user selects a small finite set of values to explore. The grid search algorithm then trains a model for every joint specification of hyperparameter values in the Cartesian product of the set of values for each individual hyperparameter. The experiment that yields the best validation set

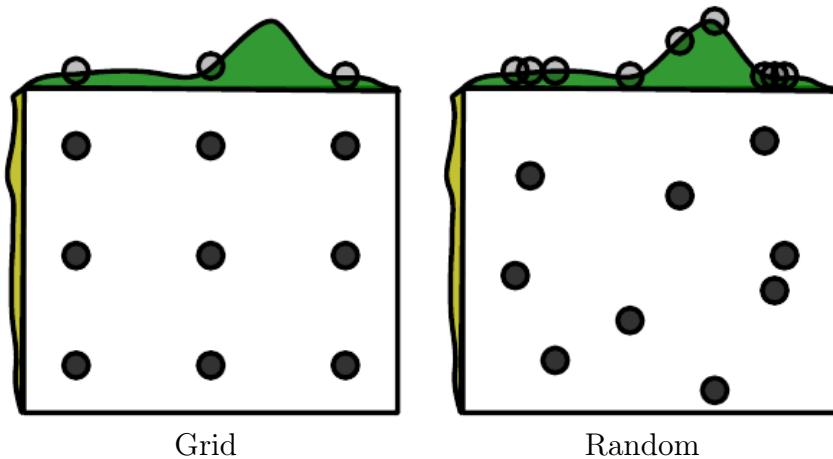


Figure 11.2: Comparison of grid search and random search. For illustration purposes we display two hyperparameters but we are typically interested in having many more. (*Left*) To perform grid search, we provide a set of values for each hyperparameter. The search algorithm runs training for every joint hyperparameter setting in the cross product of these sets. (*Right*) To perform random search, we provide a probability distribution over joint hyperparameter configurations. Usually most of these hyperparameters are independent from each other. Common choices for the distribution over a single hyperparameter include uniform and log-uniform (to sample from a log-uniform distribution, take the exp of a sample from a uniform distribution). The search algorithm then randomly samples joint hyperparameter configurations and runs training with each of them. Both grid search and random search evaluate the validation set error and return the best configuration. The figure illustrates the typical case where only some hyperparameters have a significant influence on the result. In this illustration, only the hyperparameter on the horizontal axis has a significant effect. Grid search wastes an amount of computation that is exponential in the number of non-influential hyperparameters, while random search tests a unique value of every influential hyperparameter on nearly every trial.

error is then chosen as having found the best hyperparameters. See the left of Fig. 11.2 for an illustration of a grid of hyperparameter values.

How should the lists of values to search over be chosen? In the case of numerical (ordered) hyperparameters, the smallest and largest element of each list is chosen conservatively, based on prior experience with similar experiments, to make sure that the optimal value is very likely to be in the selected range. Typically, a grid search involves picking values approximately on a **logarithmic scale**, e.g., a learning rate taken within the set  $\{.1, .01, 10^{-3}, 10^{-4}, 10^{-5}\}$ , or a number of hidden units taken with the set  $\{50, 100, 200, 500, 1000, 2000\}$ .

Grid search usually performs best when it is performed repeatedly. For example, suppose that we ran a grid search over a hyperparameter  $\alpha$  using values of  $\{-1, 0, 1\}$ . If the best value found is 1, then we underestimated the range in which the best  $\alpha$  lies and we should shift the grid and run another search with  $\alpha$  in, for example,  $\{1, 2, 3\}$ . If we find that the best value of  $\alpha$  is 0, then we may wish to refine our estimate by zooming in and running a grid search over  $\{-1, 0, 1\}$ .

The obvious problem with grid search is that its computational cost grows exponentially with the number of hyperparameters. If there are  $m$  hyperparameters, each taking at most  $n$  values, then the number of training and evaluation trials required grows as  $O(n^m)$ . The trials may be run in parallel and exploit loose parallelism (with almost no need for communication between different machines carrying out the search). Unfortunately, due to the exponential cost of grid search, even parallelization may not provide a satisfactory size of search.

#### 11.4.4 Random Search

Fortunately, there is an alternative to grid search that is as simple to program, more convenient to use, and converges much faster to good values of the hyperparameters: random search (Bergstra and Bengio, 2012).

A random search proceeds as follows. First we define a marginal distribution for each hyperparameter, e.g., a Bernoulli or multinoulli for binary or discrete hyperparameters, or a uniform distribution on a log-scale for positive real-valued hyperparameters. For example,

$$\text{log\_learning\_rate} \sim u(-1, -5) \tag{11.2}$$

$$\text{learning\_rate} = 10^{\text{log\_learning\_rate}}. \tag{11.3}$$

where  $u(a, b)$  indicates a sample of the uniform distribution in the interval  $(a, b)$ . Similarly the `log_number_of_hidden_units` may be sampled from  $u(\log(50), \log(2000))$ .

Unlike in the case of a grid search, one **should not discretize** or bin the values of the hyperparameters. This allows one to explore a larger set of values, and does not incur additional computational cost. In fact, as illustrated in Fig. 11.2, a random search can be exponentially more efficient than a grid search, when there are several hyperparameters that do not strongly affect the performance measure. This is studied at length in [Bergstra and Bengio \(2012\)](#), who found that random search reduces the validation set error much faster than grid search, in terms of the number of trials run by each method.

As with grid search, one may often want to run repeated versions of random search, to refine the search based on the results of the first run.

The main reason why random search finds good solutions faster than grid search is that there are no wasted experimental runs, unlike in the case of grid search, when two values of a hyperparameter (given values of the other hyperparameters) would give the same result. In the case of grid search, the other hyperparameters would have the same values for these two runs, whereas with random search, they would usually have different values. Hence if the change between these two values does not marginally make much difference in terms of validation set error, grid search will unnecessarily repeat two equivalent experiments while random search will still give two independent explorations of the other hyperparameters.

#### 11.4.5 Model-Based Hyperparameter Optimization

The search for good hyperparameters can be cast as an optimization problem. The decision variables are the hyperparameters. The cost to be optimized is the validation set error that results from training using these hyperparameters. In simplified settings where it is feasible to compute the gradient of some differentiable error measure on the validation set with respect to the hyperparameters, we can simply follow this gradient ([Bengio et al., 1999; Bengio, 2000; Maclaurin et al., 2015](#)). Unfortunately, in most practical settings, this gradient is unavailable, either due to its high computation and memory cost, or due to hyperparameters having intrinsically non-differentiable interactions with the validation set error, as in the case of discrete-valued hyperparameters.

To compensate for this lack of a gradient, we can build a model of the validation set error, then propose new hyperparameter guesses by performing optimization within this model. Most model-based algorithms for hyperparameter search use a Bayesian regression model to estimate both the expected value of the validation set error for each hyperparameter and the uncertainty around this expectation. Optimization thus involves a tradeoff between exploration (proposing hyperparameters

for which there is high uncertainty, which may lead to a large improvement but may also perform poorly) and exploitation (proposing hyperparameters which the model is confident will perform as well as any hyperparameters it has seen so far—usually hyperparameters that are very similar to ones it has seen before). Contemporary approaches to hyperparameter optimization include Spearmint ([Snoek et al., 2012](#)), TPE ([Bergstra et al., 2011](#)) and SMAC ([Hutter et al., 2011](#)).

Currently, we cannot unambiguously recommend Bayesian hyperparameter optimization as an established tool for achieving better deep learning results or for obtaining those results with less effort. Bayesian hyperparameter optimization sometimes performs comparably to human experts, sometimes better, but fails catastrophically on other problems. It may be worth trying to see if it works on a particular problem but is not yet sufficiently mature or reliable. That being said, hyperparameter optimization is an important field of research that, while often driven primarily by the needs of deep learning, holds the potential to benefit not only the entire field of machine learning but the discipline of engineering in general.

One drawback common to most hyperparameter optimization algorithms with more sophistication than random search is that they require for a training experiment to run to completion before they are able to extract any information from the experiment. This is much less efficient, in the sense of how much information can be gleaned early in an experiment, than manual search by a human practitioner, since one can usually tell early on if some set of hyperparameters is completely pathological. [Swersky et al. \(2014\)](#) have introduced an early version of an algorithm that maintains a set of multiple experiments. At various time points, the hyperparameter optimization algorithm can choose to begin a new experiment, to “freeze” a running experiment that is not promising, or to “thaw” and resume an experiment that was earlier frozen but now appears promising given more information.

## 11.5 Debugging Strategies

When a machine learning system performs poorly, it is usually difficult to tell whether the poor performance is intrinsic to the algorithm itself or whether there is a bug in the implementation of the algorithm. Machine learning systems are difficult to debug for a variety of reasons.

In most cases, we do not know *a priori* what the intended behavior of the algorithm is. In fact, the entire point of using machine learning is that it will discover useful behavior that we were not able to specify ourselves. If we train a

neural network on a *new* classification task and it achieves 5% test error, we have no straightforward way of knowing if this is the expected behavior or sub-optimal behavior.

A further difficulty is that most machine learning models have multiple parts that are each adaptive. If one part is broken, the other parts can adapt and still achieve roughly acceptable performance. For example, suppose that we are training a neural net with several layers parametrized by weights  $\mathbf{W}$  and biases  $\mathbf{b}$ . Suppose further that we have manually implemented the gradient descent rule for each parameter separately, and we made an error in the update for the biases:

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha \quad (11.4)$$

where  $\alpha$  is the learning rate. This erroneous update does not use the gradient at all. It causes the biases to constantly become negative throughout learning, which is clearly not a correct implementation of any reasonable learning algorithm. The bug may not be apparent just from examining the output of the model though. Depending on the distribution of the input, the weights may be able to adapt to compensate for the negative biases.

Most debugging strategies for neural nets are designed to get around one or both of these two difficulties. Either we design a case that is so simple that the correct behavior actually can be predicted, or we design a test that exercises one part of the neural net implementation in isolation.

Some important debugging tests include:

*Visualize the model in action* : When training a model to detect objects in images, view some images with the detections proposed by the model displayed superimposed on the image. When training a generative model of speech, listen to some of the speech samples it produces. This may seem obvious, but it is easy to fall into the practice of only looking at quantitative performance measurements like accuracy or log-likelihood. Directly observing the machine learning model performing its task will help you to determine whether the quantitative performance numbers it achieves seem reasonable. Evaluation bugs can be some of the most devastating bugs because they can mislead you into believing your system is performing well when it is not.

*Visualize the worst mistakes* : Most models are able to output some sort of confidence measure for the task they perform. For example, classifiers based on a softmax output layer assign a probability to each class. The probability assigned to the most likely class thus gives an estimate of the confidence the model has in its classification decision. Typically, maximum likelihood training results in these values being overestimates rather than accurate probabilities of correct prediction,

but they are somewhat useful in the sense that examples that are actually less likely to be correctly labeled receive smaller probabilities under the model. By viewing the training set examples that are the hardest to model correctly, one can often discover problems with the way the data has been preprocessed or labeled. For example, the Street View transcription system originally had a problem where the address number detection system would crop the image too tightly and omit some of the digits. The transcription network then assigned very low probability to the correct answer on these images. Sorting the images to identify the most confident mistakes showed that there was a systematic problem with the cropping. Modifying the detection system to crop much wider images resulted in much better performance of the overall system, even though the transcription network needed to be able to process greater variation in the position and scale of the address numbers.

*Reasoning about software using train and test error:* It is often difficult to determine whether the underlying software is correctly implemented. Some clues can be obtained from the train and test error. If training error is low but test error is high, then it is likely that the training procedure works correctly, and the model is overfitting for fundamental algorithmic reasons. An alternative possibility is that the test error is measured incorrectly due to a problem with saving the model after training then reloading it for test set evaluation, or if the test data was prepared differently from the training data. If both train and test error are high, then it is difficult to determine whether there is a software defect or whether the model is underfitting due to fundamental algorithmic reasons. This scenario requires further tests, described next.

*Fit a tiny dataset:* If you have high error on the training set, determine whether it is due to genuine underfitting or due to a software defect. Usually even small models can be guaranteed to be able fit a sufficiently small dataset. For example, a classification dataset with only one example can be fit just by setting the biases of the output layer correctly. Usually if you cannot train a classifier to correctly label a single example, an autoencoder to successfully reproduce a single example with high fidelity, or a generative model to consistently emit samples resembling a single example, there is a software defect preventing successful optimization on the training set. This test can be extended to a small dataset with few examples.

*Compare back-propagated derivatives to numerical derivatives:* If you are using a software framework that requires you to implement your own gradient computations, or if you are adding a new operation to a differentiation library and must define its `bprop` method, then a common source of error is implementing this gradient expression incorrectly. One way to verify that these derivatives are correct

is to compare the derivatives computed by your implementation of automatic differentiation to the derivatives computed by a *finite differences*. Because

$$f'(x) = \lim_{\epsilon \rightarrow 0} \frac{f(x + \epsilon) - f(x)}{\epsilon}, \quad (11.5)$$

we can approximate the derivative by using a small, finite  $\epsilon$ :

$$f'(x) \approx \frac{f(x + \epsilon) - f(x)}{\epsilon}. \quad (11.6)$$

We can improve the accuracy of the approximation by using the *centered difference*:

$$f'(x) \approx \frac{f(x + \frac{1}{2}\epsilon) - f(x - \frac{1}{2}\epsilon)}{\epsilon}. \quad (11.7)$$

The perturbation size  $\epsilon$  must chosen to be large enough to ensure that the perturbation is not rounded down too much by finite-precision numerical computations.

Usually, we will want to test the gradient or Jacobian of a vector-valued function  $g : \mathbb{R}^m \rightarrow \mathbb{R}^n$ . Unfortunately, finite differencing only allows us to take a single derivative at a time. We can either run finite differencing  $mn$  times to evaluate all of the partial derivatives of  $g$ , or we can apply the test to a new function that uses random projections at both the input and output of  $g$ . For example, we can apply our test of the implementation of the derivatives to  $f(x)$  where  $f(x) = \mathbf{u}^T g(\mathbf{v}x)$ , where  $\mathbf{u}$  and  $\mathbf{v}$  are randomly chosen vectors. Computing  $f'(x)$  correctly requires being able to back-propagate through  $g$  correctly, yet is efficient to do with finite differences because  $f$  has only a single input and a single output. It is usually a good idea to repeat this test for more than one value of  $\mathbf{u}$  and  $\mathbf{v}$  to reduce the chance that the test overlooks mistakes that are orthogonal to the random projection.

If one has access to numerical computation on complex numbers, then there is a very efficient way to numerically estimate the gradient by using complex numbers as input to the function ([Squire and Trapp, 1998](#)). The method is based on the observation that

$$f(x + i\epsilon) = f(x) + i\epsilon f'(x) + O(\epsilon^2) \quad (11.8)$$

$$\text{real}(f(x + i\epsilon)) = f(x) + O(\epsilon^2), \quad \text{imag}\left(\frac{f(x + i\epsilon)}{\epsilon}\right) = f'(x) + O(\epsilon^2), \quad (11.9)$$

where  $i = \sqrt{-1}$ . Unlike in the real-valued case above, there is no cancellation effect due to taking the difference between the value of  $f$  at different points. This allows the use of tiny values of  $\epsilon$  like  $\epsilon = 10^{-150}$ , which make the  $O(\epsilon^2)$  error insignificant for all practical purposes.

*Monitor histograms of activations and gradient:* It is often useful to visualize statistics of neural network activations and gradients, collected over a large amount of training iterations (maybe one epoch). The pre-activation value of hidden units can tell us if the units saturate, or how often they do. For example, for rectifiers, how often are they off? Are there units that are always off? For tanh units, the average of the absolute value of the pre-activations tells us how saturated the unit is. In a deep network where the propagated gradients quickly grow or quickly vanish, optimization may be hampered. Finally, it is useful to compare the magnitude of parameter gradients to the magnitude of the parameters themselves. As suggested by [Bottou \(2015\)](#), we would like the magnitude of parameter updates over a minibatch to represent something like 1% of the magnitude of the parameter, not 50% or 0.001% (which would make the parameters move too slowly). It may be that some groups of parameters are moving at a good pace while others are stalled. When the data is sparse (like in natural language), some parameters may be very rarely updated, and this should be kept in mind when monitoring their evolution.

Finally, many deep learning algorithms provide some sort of guarantee about the results produced at each step. For example, in Part III, we will see some approximate inference algorithms that work by using algebraic solutions to optimization problems. Typically these can be debugged by testing each of their guarantees. Some guarantees that some optimization algorithms offer include that the objective function will never increase after one step of the algorithm, that the gradient with respect to some subset of variables will be zero after each step of the algorithm, and that the gradient with respect to all variables will be zero at convergence. Usually due to rounding error, these conditions will not hold exactly in a digital computer, so the debugging test should include some tolerance parameter.

## 11.6 Example: Multi-Digit Number Recognition

To provide an end-to-end description of how to apply our design methodology in practice, we present a brief account of the Street View transcription system, from the point of view of designing the deep learning components. Obviously, many other components of the complete system, such as the Street View cars, the database infrastructure, and so on, were of paramount importance.

From the point of view of the machine learning task, the process began with data collection. The cars collected the raw data and human operators provided labels. The transcription task was preceded by a significant amount of dataset

curation, including using other machine learning techniques to **detect** the house numbers prior to transcribing them.

The transcription project began with a choice of performance metrics and desired values for these metrics. An important general principle is to tailor the choice of metric to the business goals for the project. Because maps are only useful if they have high accuracy, it was important to set a high accuracy requirement for this project. Specifically, the goal was to obtain human-level, 98% accuracy. This level of accuracy may not always be feasible to obtain. In order to reach this level of accuracy, the Street View transcription system sacrifices coverage. Coverage thus became the main performance metric optimized during the project, with accuracy held at 98%. As the convolutional network improved, it became possible to reduce the confidence threshold below which the network refuses to transcribe the input, eventually exceeding the goal of 95% coverage.

After choosing quantitative goals, the next step in our recommended methodology is to rapidly establish a sensible baseline system. For vision tasks, this means a convolutional network with rectified linear units. The transcription project began with such a model. At the time, it was not common for a convolutional network to output a sequence of predictions. In order to begin with the simplest possible baseline, the first implementation of the output layer of the model consisted of  $n$  different softmax units to predict a sequence of  $n$  characters. These softmax units were trained exactly the same as if the task were classification, with each softmax unit trained independently.

Our recommended methodology is to iteratively refine the baseline and test whether each change makes an improvement. The first change to the Street View transcription system was motivated by a theoretical understanding of the coverage metric and the structure of the data. Specifically, the network refuses to classify an input  $\mathbf{x}$  whenever the probability of the output sequence  $p(\mathbf{y} \mid \mathbf{x}) < t$  for some threshold  $t$ . Initially, the definition of  $p(\mathbf{y} \mid \mathbf{x})$  was ad-hoc, based on simply multiplying all of the softmax outputs together. This motivated the development of a specialized output layer and cost function that actually computed a principled log-likelihood. This approach allowed the example rejection mechanism to function much more effectively.

At this point, coverage was still below 90%, yet there were no obvious theoretical problems with the approach. Our methodology therefore suggests to instrument the train and test set performance in order to determine whether the problem is underfitting or overfitting. In this case, train and test set error were nearly identical. Indeed, the main reason this project proceeded so smoothly was the availability of a dataset with tens of millions of labeled examples. Because train

and test set error were so similar, this suggested that the problem was either due to underfitting or due to a problem with the training data. One of the debugging strategies we recommend is to visualize the model’s worst errors. In this case, that meant visualizing the incorrect training set transcriptions that the model gave the highest confidence. These proved to mostly consist of examples where the input image had been cropped too tightly, with some of the digits of the address being removed by the cropping operation. For example, a photo of an address “1849” might be cropped too tightly, with only the “849” remaining visible. This problem could have been resolved by spending weeks improving the accuracy of the address number detection system responsible for determining the cropping regions. Instead, the team took a much more practical decision, to simply expand the width of the crop region to be systematically wider than the address number detection system predicted. This single change added ten percentage points to the transcription system’s coverage.

Finally, the last few percentage points of performance came from adjusting hyperparameters. This mostly consisted of making the model larger while maintaining some restrictions on its computational cost. Because train and test error remained roughly equal, it was always clear that any performance deficits were due to underfitting, as well as due to a few remaining problems with the dataset itself.

Overall, the transcription project was a great success, and allowed hundreds of millions of addresses to be transcribed both faster and at lower cost than would have been possible via human effort.

We hope that the design principles described in this chapter will lead to many other similar successes.

# Chapter 12

# Applications

In this chapter, we describe how to use deep learning to solve applications in computer vision, speech recognition, natural language processing, and other application areas of commercial interest. We begin by discussing the large scale neural network implementations required for most serious AI applications. Next, we review several specific application areas that deep learning has been used to solve. While one goal of deep learning is to design algorithms that are capable of solving a broad variety of tasks, so far some degree of specialization is needed. For example, vision tasks require processing a large number of input features (pixels) per example. Language tasks require modeling a large number of possible values (words in the vocabulary) per input feature.

## 12.1 Large Scale Deep Learning

Deep learning is based on the philosophy of connectionism: while an individual biological neuron or an individual feature in a machine learning model is not intelligent, a large population of these neurons or features acting together can exhibit intelligent behavior. It truly is important to emphasize the fact that the number of neurons must be *large*. One of the key factors responsible for the improvement in neural network's accuracy and the improvement of the complexity of tasks they can solve between the 1980s and today is the dramatic increase in the size of the networks we use. As we saw in Sec. 1.2.3, network sizes have grown exponentially for the past three decades, yet artificial neural networks are only as large as the nervous systems of insects.

Because the size of neural networks is of paramount importance, deep learning

requires high performance hardware and software infrastructure.

### 12.1.1 Fast CPU Implementations

Traditionally, neural networks were trained using the CPU of a single machine. Today, this approach is generally considered insufficient. We now mostly use GPU computing or the CPUs of many machines networked together. Before moving to these expensive setups, researchers worked hard to demonstrate that CPUs could not manage the high computational workload required by neural networks.

A description of how to implement efficient numerical CPU code is beyond the scope of this book, but we emphasize here that careful implementation for specific CPU families can yield large improvements. For example, in 2011, the best CPUs available could run neural network workloads faster when using fixed-point arithmetic rather than floating-point arithmetic. By creating a carefully tuned fixed-point implementation, [Vanhoucke et al. \(2011\)](#) obtained a  $3\times$  speedup over a strong floating-point system. Each new model of CPU has different performance characteristics, so sometimes floating-point implementations can be faster too. The important principle is that careful specialization of numerical computation routines can yield a large payoff. Other strategies, besides choosing whether to use fixed or floating point, include optimizing data structures to avoid cache misses and using vector instructions. Many machine learning researchers neglect these implementation details, but when the performance of an implementation restricts the size of the model, the accuracy of the model suffers.

### 12.1.2 GPU Implementations

Most modern neural network implementations are based on graphics processing units. Graphics processing units (GPUs) are specialized hardware components that were originally developed for graphics applications. The consumer market for video gaming systems spurred development of graphics processing hardware. The performance characteristics needed for good video gaming systems turn out to be beneficial for neural networks as well.

Video game rendering requires performing many operations in parallel quickly. Models of characters and environments are specified in terms of lists of 3-D coordinates of vertices. Graphics cards must perform matrix multiplication and division on many vertices in parallel to convert these 3-D coordinates into 2-D on-screen coordinates. The graphics card must then perform many computations at each pixel in parallel to determine the color of each pixel. In both cases, the

computations are fairly simple and do not involve much branching compared to the computational workload that a CPU usually encounters. For example, each vertex in the same rigid object will be multiplied by the same matrix; there is no need to evaluate an if statement per-vertex to determine which matrix to multiply by. The computations are also entirely independent of each other, and thus may be parallelized easily. The computations also involve processing massive buffers of memory, containing bitmaps describing the texture (color pattern) of each object to be rendered. Together, this results in graphics cards having been designed to have a high degree of parallelism and high memory bandwidth, at the cost of having a lower clock speed and less branching capability relative to traditional CPUs.

Neural network algorithms require the same performance characteristics as the real-time graphics algorithms described above. Neural networks usually involve large and numerous buffers of parameters, activation values, and gradient values, each of which must be completely updated during every step of training. These buffers are large enough to fall outside the cache of a traditional desktop computer so the memory bandwidth of the system often becomes the rate limiting factor. GPUs offer a compelling advantage over CPUs due to their high memory bandwidth. Neural network training algorithms typically do not involve much branching or sophisticated control, so they are appropriate for GPU hardware. Since neural networks can be divided into multiple individual “neurons” that can be processed independently from the other neurons in the same layer, neural networks easily benefit from the parallelism of GPU computing.

GPU hardware was originally so specialized that it could only be used for graphics tasks. Over time, GPU hardware became more flexible, allowing custom subroutines to be used to transform the coordinates of vertices or assign colors to pixels. In principle, there was no requirement that these pixel values actually be based on a rendering task. These GPUs could be used for scientific computing by writing the output of a computation to a buffer of pixel values. [Steinkrau et al. \(2005\)](#) implemented a two-layer fully connected neural network on a GPU and reported a 3X speedup over their CPU-based baseline. Shortly thereafter, [Chellapilla et al. \(2006\)](#) demonstrated that the same technique could be used to accelerate supervised convolutional networks.

The popularity of graphics cards for neural network training exploded after the advent of *general purpose GPUs*. These GP-GPUs could execute arbitrary code, not just rendering subroutines. NVIDIA’s CUDA programming language provided a way to write this arbitrary code in a C-like language. With their relatively convenient programming model, massive parallelism, and high memory bandwidth,

GP-GPUs now offer an ideal platform for neural network programming. This platform was rapidly adopted by deep learning researchers soon after it became available (Raina *et al.*, 2009; Ciresan *et al.*, 2010).

Writing efficient code for GP-GPUs remains a difficult task best left to specialists. The techniques required to obtain good performance on GPU are very different from those used on CPU. For example, good CPU-based code is usually designed to read information from the cache as much as possible. On GPU, most writable memory locations are not cached, so it can actually be faster to compute the same value twice, rather than compute it once and read it back from memory. GPU code is also inherently multi-threaded and the different threads must be coordinated with each other carefully. For example, memory operations are faster if they can be *coalesced*. Coalesced reads or writes occur when several threads can each read or write a value that they need simultaneously, as part of a single memory transaction. Different models of GPUs are able to coalesce different kinds of read or write patterns. Typically, memory operations are easier to coalesce if among  $n$  threads, thread  $i$  accesses byte  $i + j$  of memory, and  $j$  is a multiple of some power of 2. The exact specifications differ between models of GPU. Another common consideration for GPUs is making sure that each thread in a group executes the same instruction simultaneously. This means that branching can be difficult on GPU. Threads are divided into small groups called *warps*. Each thread in a warp executes the same instruction during each cycle, so if different threads within the same warp need to execute different code paths, these different code paths must be traversed sequentially rather than in parallel.

Due to the difficulty of writing high performance GPU code, researchers should structure their workflow to avoid needing to write new GPU code in order to test new models or algorithms. Typically, one can do this by building a software library of high performance operations like convolution and matrix multiplication, then specifying models in terms of calls to this library of operations. For example, the machine learning library Pylearn2 (Goodfellow *et al.*, 2013c) specifies all of its machine learning algorithms in terms of calls to Theano (Bergstra *et al.*, 2010; Bastien *et al.*, 2012) and cuda-convnet (Krizhevsky, 2010), which provide these high-performance operations. This factored approach can also ease support for multiple kinds of hardware. For example, the same Theano program can run on either CPU or GPU, without needing to change any of the calls to Theano itself. Other libraries like TensorFlow (Abadi *et al.*, 2015) and Torch (Collobert *et al.*, 2011b) provide similar features.

### 12.1.3 Large Scale Distributed Implementations

In many cases, the computational resources available on a single machine are insufficient. We therefore want to distribute the workload of training and inference across many machines.

Distributing inference is simple, because each input example we want to process can be run by a separate machine. This is known as *data parallelism*.

It is also possible to get *model parallelism*, where multiple machines work together on a single datapoint, with each machine running a different part of the model. This is feasible for both inference and training.

Data parallelism during training is somewhat harder. We can increase the size of the minibatch used for a single SGD step, but usually we get less than linear returns in terms of optimization performance. It would be better to allow multiple machines to compute multiple gradient descent steps in parallel. Unfortunately, the standard definition of gradient descent is as a completely sequential algorithm: the gradient at step  $t$  is a function of the parameters produced by step  $t - 1$ .

This can be solved using *asynchronous stochastic gradient descent* (Bengio *et al.*, 2001; Recht *et al.*, 2011). In this approach, several processor cores share the memory representing the parameters. Each core reads parameters without a lock, then computes a gradient, then increments the parameters without a lock. This reduces the average amount of improvement that each gradient descent step yields, because some of the cores overwrite each other's progress, but the increased rate of production of steps causes the learning process to be faster overall. Dean *et al.* (2012) pioneered the multi-machine implementation of this lock-free approach to gradient descent, where the parameters are managed by a *parameter server* rather than stored in shared memory. Distributed asynchronous gradient descent remains the primary strategy for training large deep networks and is used by most major deep learning groups in industry (Chilimbi *et al.*, 2014; Wu *et al.*, 2015). Academic deep learning researchers typically cannot afford the same scale of distributed learning systems but some research has focused on how to build distributed networks with relatively low-cost hardware available in the university setting (Coates *et al.*, 2013).

### 12.1.4 Model Compression

In many commercial applications, it is much more important that the time and memory cost of running inference in a machine learning model be low than that the time and memory cost of training be low. For applications that do not require

personalization, it is possible to train a model once, then deploy it to be used by billions of users. In many cases, the end user is more resource-constrained than the developer. For example, one might train a speech recognition network with a powerful computer cluster, then deploy it on mobile phones.

A key strategy for reducing the cost of inference is *model compression* (Buciluă *et al.*, 2006). The basic idea of model compression is to replace the original, expensive model with a smaller model that requires less memory and runtime to store and evaluate.

Model compression is applicable when the size of the original model is driven primarily by a need to prevent overfitting. In most cases, the model with the lowest generalization error is an ensemble of several independently trained models. Evaluating all  $n$  ensemble members is expensive. Sometimes, even a single model generalizes better if it is large (for example, if it is regularized with dropout).

These large models learn some function  $f(\mathbf{x})$ , but do so using many more parameters than are necessary for the task. Their size is necessary only due to the limited number of training examples. As soon as we have fit this function  $f(\mathbf{x})$ , we can generate a training set containing infinitely many examples, simply by applying  $f$  to randomly sampled points  $\mathbf{x}$ . We then train the new, smaller, model to match  $f(\mathbf{x})$  on these points. In order to most efficiently use the capacity of the new, small model, it is best to sample the new  $\mathbf{x}$  points from a distribution resembling the actual test inputs that will be supplied to the model later. This can be done by corrupting training examples or by drawing points from a generative model trained on the original training set.

Alternatively, one can train the smaller model only on the original training points, but train it to copy other features of the model, such as its posterior distribution over the incorrect classes (Hinton *et al.*, 2014, 2015).

### 12.1.5 Dynamic Structure

One strategy for accelerating data processing systems in general is to build systems that have *dynamic structure* in the graph describing the computation needed to process an input. Data processing systems can dynamically determine which subset of many neural networks should be run on a given input. Individual neural networks can also exhibit dynamic structure internally by determining which subset of features (hidden units) to compute given information from the input. This form of dynamic structure inside neural networks is sometimes called *conditional computation* (Bengio, 2013; Bengio *et al.*, 2013b). Since many components of the architecture may be relevant only for a small amount of possible inputs, the system

can run faster by computing these features only when they are needed.

Dynamic structure of computations is a basic computer science principle applied generally throughout the software engineering discipline. The simplest versions of dynamic structure applied to neural networks are based on determining which subset of some group of neural networks (or other machine learning models) should be applied to a particular input.

A venerable strategy for accelerating inference in a classifier is to use a *cascade* of classifiers. The cascade strategy may be applied when the goal is to detect the presence of a rare object (or event). To know for sure that the object is present, we must use a sophisticated classifier with high capacity, that is expensive to run. However, because the object is rare, we can usually use much less computation to reject inputs as not containing the object. In these situations, we can train a sequence of classifiers. The first classifiers in the sequence have low capacity, and are trained to have high recall. In other words, they are trained to make sure we do not wrongly reject an input when the object is present. The final classifier is trained to have high precision. At test time, we run inference by running the classifiers in a sequence, abandoning any example as soon as any one element in the cascade rejects it. Overall, this allows us to verify the presence of objects with high confidence, using a high capacity model, but does not force us to pay the cost of full inference for every example. There are two different ways that the cascade can achieve high capacity. One way is to make the later members of the cascade individually have high capacity. In this case, the system as a whole obviously has high capacity, because some of its individual members do. It is also possible to make a cascade in which every individual model has low capacity but the system as a whole has high capacity due to the combination of many small models. [Viola and Jones \(2001\)](#) used a cascade of boosted decision trees to implement a fast and robust face detector suitable for use in handheld digital cameras. Their classifier localizes a face using essentially a sliding window approach in which many windows are examined and rejected if they do not contain faces. Another version of cascades uses the earlier models to implement a sort of hard attention mechanism: the early members of the cascade localize an object and later members of the cascade perform further processing given the location of the object. For example, Google transcribes address numbers from Street View imagery using a two-step cascade that first locates the address number with one machine learning model and then transcribes it with another ([Goodfellow et al., 2014d](#)).

Decision trees themselves are an example of dynamic structure, because each node in the tree determines which of its subtrees should be evaluated for each input. A simple way to accomplish the union of deep learning and dynamic structure

is to train a decision tree in which each node uses a neural network to make the splitting decision (Guo and Gelfand, 1992), though this has typically not been done with the primary goal of accelerating inference computations.

In the same spirit, one can use a neural network, called the *gater* to select which one out of several *expert networks* will be used to compute the output, given the current input. The first version of this idea is called the *mixture of experts* (Nowlan, 1990; Jacobs *et al.*, 1991), in which the gater outputs a set of probabilities or weights (obtained via a softmax nonlinearity), one per expert, and the final output is obtained by the weighted combination of the output of the experts. In that case, the use of the gater does not offer a reduction in computational cost, but if a single expert is chosen by the gater for each example, we obtain the *hard mixture of experts* (Collobert *et al.*, 2001, 2002), which can considerably accelerate training and inference time. This strategy works well when the number of gating decisions is small because it is not combinatorial. But when we want to select different subsets of units or parameters, it is not possible to use a “soft switch” because it requires enumerating (and computing outputs for) all the gater configurations. To deal with this problem, several approaches have been explored to train combinatorial gaters. Bengio *et al.* (2013b) experiment with several estimators of the gradient on the gating probabilities, while Bacon *et al.* (2015) and Bengio *et al.* (2015a) use reinforcement learning techniques (policy gradient) to learn a form of conditional dropout on blocks of hidden units and get an actual reduction in computational cost without impacting negatively on the quality of the approximation.

Another kind of dynamic structure is a switch, where a hidden unit can receive input from different units depending on the context. This dynamic routing approach can be interpreted as an attention mechanism (Olshausen *et al.*, 1993). So far, the use of a hard switch has not proven effective on large-scale applications. Contemporary approaches instead use a weighted average over many possible inputs, and thus do not achieve all of the possible computational benefits of dynamic structure. Contemporary attention mechanisms are described in Sec. 12.4.5.1.

One major obstacle to using dynamically structured systems is the decreased degree of parallelism that results from the system following different code branches for different inputs. This means that few operations in the network can be described as matrix multiplication or batch convolution on a minibatch of examples. We can write more specialized sub-routines that convolve each example with different kernels or multiply each row of a design matrix by a different set of columns of weights. Unfortunately, these more specialized subroutines are difficult to implement efficiently. CPU implementations will be slow due to the lack of cache coherence and GPU implementations will be slow due to the lack of coalesced

memory transactions and the need to serialize warps when members of a warp take different branches. In some cases, these issues can be mitigated by partitioning the examples into groups that all take the same branch, and processing these groups of examples simultaneously. This can be an acceptable strategy for minimizing the time required to process a fixed amount of examples in an offline setting. In a real-time setting where examples must be processed continuously, partitioning the workload can result in load-balancing issues. For example, if we assign one machine to process the first step in a cascade and another machine to process the last step in a cascade, then the first will tend to be overloaded and the last will tend to be underloaded. Similar issues arise if each machine is assigned to implement different nodes of a neural decision tree.

### 12.1.6 Specialized Hardware Implementations of Deep Networks

Since the early days of neural networks research, hardware designers have worked on specialized hardware implementations that could speed up training and/or inference of neural network algorithms. See early and more recent reviews of specialized hardware for deep networks ([Lindsey and Lindblad, 1994](#); [Beiu \*et al.\*, 2003](#); [Misra and Saha, 2010](#)).

Different forms of specialized hardware ([Graf and Jackel, 1989](#); [Mead and Ismail, 2012](#); [Kim \*et al.\*, 2009](#); [Pham \*et al.\*, 2012](#); [Chen \*et al.\*, 2014a,b](#)) have been developed over the last decades, either with ASICs (application-specific integrated circuit), either with digital (based on binary representations of numbers), analog ([Graf and Jackel, 1989](#); [Mead and Ismail, 2012](#)) (based on physical implementations of continuous values as voltages or currents) or hybrid implementations (combining digital and analog components). In recent years more flexible FPGA (field programmable gated array) implementations (where the particulars of the circuit can be written on the chip after it has been built) have been developed.

Though software implementations on general-purpose processing units (CPUs and GPUs) typically use 32 or 64 bits of precision to represent floating point numbers, it has long been known that it was possible to use less precision, at least at inference time ([Holt and Baker, 1991](#); [Holi and Hwang, 1993](#); [Presley and Haggard, 1994](#); [Simard and Graf, 1994](#); [Wawrynek \*et al.\*, 1996](#); [Savich \*et al.\*, 2007](#)). This has become a more pressing issue in recent years as deep learning has gained in popularity in industrial products, and as the great impact of faster hardware was demonstrated with GPUs. Another factor that motivates current research on specialized hardware for deep networks is that the rate of progress of a single CPU or GPU core has slowed down, and most recent improvements in computing speed have come from parallelization across cores (either in CPUs or

GPUs). This is very different from the situation of the 1990s (the previous neural network era) where the hardware implementations of neural networks (which might take two years from inception to availability of a chip) could not keep up with the rapid progress and low prices of general-purpose CPUs. Building specialized hardware is thus a way to push the envelope further, at a time when new hardware designs are being developed for low-power devices such as phones, aiming for general-public applications of deep learning (e.g., with speech, computer vision or natural language).

Recent work on low-precision implementations of backprop-based neural nets (Vanhoucke *et al.*, 2011; Courbariaux *et al.*, 2015; Gupta *et al.*, 2015) suggests that between 8 and 16 bits of precision can suffice for using or training deep neural networks with back-propagation. What is clear is that more precision is required during training than at inference time, and that some forms of dynamic fixed point representation of numbers can be used to reduce how many bits are required per number. Traditional fixed point numbers are restricted to a fixed range (which corresponds to a given exponent in a floating point representation). Dynamic fixed point representations share that range among a set of numbers (such as all the weights in one layer). Using fixed point rather than floating point representations and using less bits per number reduces the hardware surface area, power requirements and computing time needed for performing multiplications, and multiplications are the most demanding of the operations needed to use or train a modern deep network with backprop.

## 12.2 Computer Vision

Computer vision has traditionally been one of the most active research areas for deep learning applications, because vision is a task that is effortless for humans and many animals but challenging for computers (Ballard *et al.*, 1983). Many of the most popular standard benchmark tasks for deep learning algorithms are forms of object recognition or optical character recognition.

Computer vision is a very broad field encompassing a wide variety of ways of processing images, and an amazing diversity of applications. Applications of computer vision range from reproducing human visual abilities, such as recognizing faces, to creating entirely new categories of visual abilities. As an example of the latter category, one recent computer vision application is to recognize sound waves from the vibrations they induce in objects visible in a video (Davis *et al.*, 2014). Most deep learning research on computer vision has not focused on such exotic applications that expand the realm of what is possible with imagery but

rather a small core of AI goals aimed at replicating human abilities. Most deep learning for computer vision is used for object recognition or detection of some form, whether this means reporting which object is present in an image, annotating an image with bounding boxes around each object, transcribing a sequence of symbols from an image, or labeling each pixel in an image with the identity of the object it belongs to. Because generative modeling has been a guiding principle of deep learning research, there is also a large body of work on image synthesis using deep models. While image synthesis *ex nihilo* is usually not considered a computer vision endeavor, models capable of image synthesis are usually useful for image restoration, a computer vision task involving repairing defects in images or removing objects from images.

### 12.2.1 Preprocessing

Many application areas require sophisticated preprocessing because the original input comes in a form that is difficult for many deep learning architectures to represent. Computer vision usually requires relatively little of this kind of preprocessing. The images should be standardized so that their pixels all lie in the same, reasonable range, like  $[0,1]$  or  $[-1, 1]$ . Mixing images that lie in  $[0,1]$  with images that lie in  $[0, 255]$  will usually result in failure. Formatting images to have the same scale is the only kind of preprocessing that is strictly necessary. Many computer vision architectures require images of a standard size, so images must be cropped or scaled to fit that size. However, even this rescaling is not always strictly necessary. Some convolutional models accept variably-sized inputs and dynamically adjust the size of their pooling regions to keep the output size constant (Waibel *et al.*, 1989). Other convolutional models have variable-sized output that automatically scales in size with the input, such as models that denoise or label each pixel in an image (Hadsell *et al.*, 2007).

Dataset augmentation may be seen as a way of preprocessing the training set only. Dataset augmentation is an excellent way to reduce the generalization error of most computer vision models. A related idea applicable at test time is to show the model many different versions of the same input (for example, the same image cropped at slightly different locations) and have the different instantiations of the model vote to determine the output. This latter idea can be interpreted as an ensemble approach, and helps to reduce generalization error.

Other kinds of preprocessing are applied to both the train and the test set with the goal of putting each example into a more canonical form in order to reduce the amount of variation that the model needs to account for. Reducing the amount of variation in the data can both reduce generalization error and reduce the size of

the model needed to fit the training set. Simpler tasks may be solved by smaller models, and simpler solutions are more likely to generalize well. Preprocessing of this kind is usually designed to remove some kind of variability in the input data that is easy for a human designer to describe and that the human designer is confident has no relevance to the task. When training with large datasets and large models, this kind of preprocessing is often unnecessary, and it is best to just let the model learn which kinds of variability it should become invariant to. For example, the AlexNet system for classifying ImageNet only has one preprocessing step: subtracting the mean across training examples of each pixel ([Krizhevsky et al., 2012](#)).

### 12.2.1.1 Contrast Normalization

One of the most obvious sources of variation that can be safely removed for many tasks is the amount of contrast in the image. Contrast simply refers to the magnitude of the difference between the bright and the dark pixels in an image. There are many ways of quantifying the contrast of an image. In the context of deep learning, contrast usually refers to the standard deviation of the pixels in an image or region of an image. Suppose we have an image represented by a tensor  $\mathbf{X} \in \mathbb{R}^{r \times c \times 3}$ , with  $X_{i,j,1}$  being the red intensity at row  $i$  and column  $j$ ,  $X_{i,j,2}$  giving the green intensity and  $X_{i,j,3}$  giving the blue intensity. Then the contrast of the entire image is given by

$$\sqrt{\frac{1}{3rc} \sum_{i=1}^r \sum_{j=1}^c \sum_{k=1}^3 (X_{i,j,k} - \bar{\mathbf{X}})^2} \quad (12.1)$$

where  $\bar{\mathbf{X}}$  is the mean intensity of the entire image:

$$\bar{\mathbf{X}} = \frac{1}{3rc} \sum_{i=1}^r \sum_{j=1}^c \sum_{k=1}^3 X_{i,j,k}. \quad (12.2)$$

*Global contrast normalization* (GCN) aims to prevent images from having varying amounts of contrast by subtracting the mean from each image, then rescaling it so that the standard deviation across its pixels is equal to some constant  $s$ . This approach is complicated by the fact that no scaling factor can change the contrast of a zero-contrast image (one whose pixels all have equal intensity). Images with very low but non-zero contrast often have little information content. Dividing by the true standard deviation usually accomplishes nothing more than amplifying sensor noise or compression artifacts in such cases. This

motivates introducing a small, positive regularization parameter  $\lambda$  to bias the estimate of the standard deviation. Alternately, one can constrain the denominator to be at least  $\epsilon$ . Given an input image  $\mathbf{X}$ , GCN produces an output image  $\mathbf{X}'$ , defined such that

$$\mathbf{X}'_{i,j,k} = s \frac{\mathbf{X}_{i,j,k} - \bar{\mathbf{X}}}{\max \left\{ \epsilon, \sqrt{\lambda + \frac{1}{3rc} \sum_{i=1}^r \sum_{j=1}^c \sum_{k=1}^3 (\mathbf{X}_{i,j,k} - \bar{\mathbf{X}})^2} \right\}}. \quad (12.3)$$

Datasets consisting of large images cropped to interesting objects are unlikely to contain any images with nearly constant intensity. In these cases, it is safe to practically ignore the small denominator problem by setting  $\lambda = 0$  and avoid division by 0 in extremely rare cases by setting  $\epsilon$  to an extremely low value like  $10^{-8}$ . This is the approach used by Goodfellow *et al.* (2013a) on the CIFAR-10 dataset. Small images cropped randomly are more likely to have nearly constant intensity, making aggressive regularization more useful. Coates *et al.* (2011) used  $\epsilon = 0$  and  $\lambda = 10$  on small, randomly selected patches drawn from CIFAR-10.

The scale parameter  $s$  can usually be set to 1, as done by Coates *et al.* (2011), or chosen to make each individual pixel have standard deviation across examples close to 1, as done by Goodfellow *et al.* (2013a).

The standard deviation in Eq. 12.3 is just a rescaling of the  $L^2$  norm of the image (assuming the mean of the image has already been removed). It is preferable to define GCN in terms of standard deviation rather than  $L^2$  norm because the standard deviation includes division by the number of pixels, so GCN based on standard deviation allows the same  $s$  to be used regardless of image size. However, the observation that the  $L^2$  norm is proportional to the standard deviation can help build a useful intuition. One can understand GCN as mapping examples to a spherical shell. See Fig. 12.1 for an illustration. This can be a useful property because neural networks are often better at responding to directions in space rather than exact locations. Responding to multiple distances in the same direction requires hidden units with collinear weight vectors but different biases. Such coordination can be difficult for the learning algorithm to discover. Additionally, many shallow graphical models have problems with representing multiple separated modes along the same line. GCN avoids these problems by reducing each example to a direction rather than a direction and a distance.

Counterintuitively, there is a preprocessing operation known as *sphering* and it is not the same operation as GCN. Sphering does not refer to making the data lie on a spherical shell, but rather to rescaling the principal components to have equal variance, so that the multivariate normal distribution used by PCA has spherical contours. Sphering is more commonly known as *whitening*.

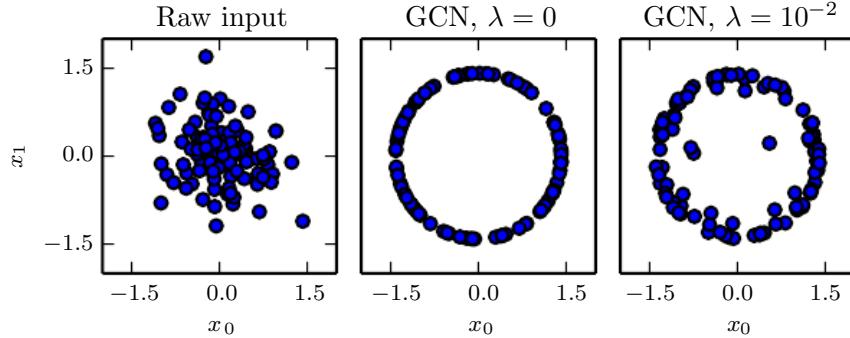


Figure 12.1: GCN maps examples onto a sphere. (*Left*) Raw input data may have any norm. (*Center*) GCN with  $\lambda = 0$  maps all non-zero examples perfectly onto a sphere. Here we use  $s = 1$  and  $\epsilon = 10^{-8}$ . Because we use GCN based on normalizing the standard deviation rather than the  $L^2$  norm, the resulting sphere is not the unit sphere. (*Right*) Regularized GCN, with  $\lambda > 0$ , draws examples toward the sphere but does not completely discard the variation in their norm. We leave  $s$  and  $\epsilon$  the same as before.

Global contrast normalization will often fail to highlight image features we would like to stand out, such as edges and corners. If we have a scene with a large dark area and a large bright area (such as a city square with half the image in the shadow of a building) then global contrast normalization will ensure there is a large difference between the brightness of the dark area and the brightness of the light area. It will not, however, ensure that edges within the dark region stand out.

This motivates *local contrast normalization*. Local contrast normalization ensures that the contrast is normalized across each small window, rather than over the image as a whole. See Fig. 12.2 for a comparison of global and local contrast normalization.

Various definitions of local contrast normalization are possible. In all cases, one modifies each pixel by subtracting a mean of nearby pixels and dividing by a standard deviation of nearby pixels. In some cases, this is literally the mean and standard deviation of all pixels in a rectangular window centered on the pixel to be modified (Pinto *et al.*, 2008). In other cases, this is a weighted mean and weighted standard deviation using Gaussian weights centered on the pixel to be modified. In the case of color images, some strategies process different color channels separately while others combine information from different channels to normalize each pixel (Sermanet *et al.*, 2012).

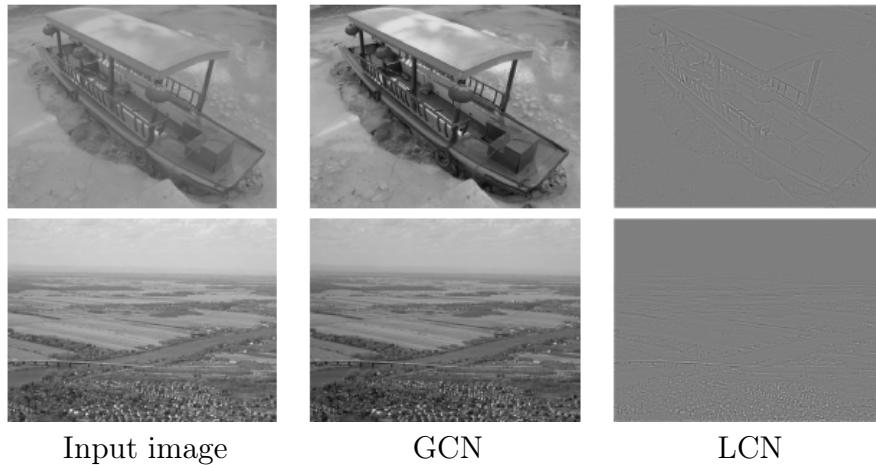


Figure 12.2: A comparison of global and local contrast normalization. Visually, the effects of global contrast normalization are subtle. It places all images on roughly the same scale, which reduces the burden on the learning algorithm to handle multiple scales. Local contrast normalization modifies the image much more, discarding all regions of constant intensity. This allows the model to focus on just the edges. Regions of fine texture, such as the houses in the second row, may lose some detail due to the bandwidth of the normalization kernel being too high.

Local contrast normalization can usually be implemented efficiently by using separable convolution (see Sec. 9.8) to compute feature maps of local means and local standard deviations, then using element-wise subtraction and element-wise division on different feature maps.

Local contrast normalization is a differentiable operation and can also be used as a nonlinearity applied to the hidden layers of a network, as well as a preprocessing operation applied to the input.

As with global contrast normalization, we typically need to regularize local contrast normalization to avoid division by zero. In fact, because local contrast normalization typically acts on smaller windows, it is even more important to regularize. Smaller windows are more likely to contain values that are all nearly the same as each other, and thus more likely to have zero standard deviation.

### 12.2.1.2 Dataset Augmentation

As described in Sec. 7.4, it is easy to improve the generalization of a classifier by increasing the size of the training set by adding extra copies of the training examples that have been modified with transformations that do not change the

class. Object recognition is a classification task that is especially amenable to this form of dataset augmentation because the class is invariant to so many transformations and the input can be easily transformed with many geometric operations. As described before, classifiers can benefit from random translations, rotations, and in some cases, flips of the input to augment the dataset. In specialized computer vision applications, more advanced transformations are commonly used for dataset augmentation. These schemes include random perturbation of the colors in an image (Krizhevsky *et al.*, 2012) and nonlinear geometric distortions of the input (LeCun *et al.*, 1998b).

## 12.3 Speech Recognition

The task of speech recognition is to map an acoustic signal containing a spoken natural language utterance into the corresponding sequence of words intended by the speaker. Let  $\mathbf{X} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(T)})$  denote the sequence of acoustic input vectors (traditionally produced by splitting the audio into 20ms frames). Most speech recognition systems preprocess the input using specialized hand-designed features, but some (Jaitly and Hinton, 2011) deep learning systems learn features from raw input. Let  $\mathbf{y} = (y_1, y_2, \dots, y_N)$  denote the target output sequence (usually a sequence of words or characters). The *automatic speech recognition* (ASR) task consists of creating a function  $f_{\text{ASR}}^*$  that computes the most probable linguistic sequence  $\mathbf{y}$  given the acoustic sequence  $\mathbf{X}$ :

$$f_{\text{ASR}}^*(\mathbf{X}) = \arg \max_{\mathbf{y}} P^*(\mathbf{y} \mid \mathbf{X} = \mathbf{X}) \quad (12.4)$$

where  $P^*$  is the true conditional distribution relating the inputs  $\mathbf{X}$  to the targets  $\mathbf{y}$ .

Since the 1980s and until about 2009–2012, state-of-the art speech recognition systems primarily combined hidden Markov models (HMMs) and Gaussian mixture models (GMMs). GMMs modeled the association between acoustic features and phonemes (Bahl *et al.*, 1987), while HMMs modeled the sequence of phonemes. The GMM-HMM model family treats acoustic waveforms as being generated by the following process: first an HMM generates a sequence of phonemes and discrete sub-phonemic states (such as the beginning, middle, and end of each phoneme), then a GMM transforms each discrete symbol into a brief segment of audio waveform. Although GMM-HMM systems dominated ASR until recently, speech recognition was actually one of the first areas where neural networks were applied, and numerous ASR systems from the late 1980s and early 1990s used

neural nets (Bourlard and Wellekens, 1989; Waibel *et al.*, 1989; Robinson and Fallside, 1991; Bengio *et al.*, 1991, 1992; Konig *et al.*, 1996). At the time, the performance of ASR based on neural nets approximately matched the performance of GMM-HMM systems. For example, Robinson and Fallside (1991) achieved 26% phoneme error rate on the TIMIT (Garofolo *et al.*, 1993) corpus (with 39 phonemes to discriminate between), which was better than or comparable to HMM-based systems. Since then, TIMIT has been a benchmark for phoneme recognition, playing a role similar to the role MNIST plays for object recognition. However, because of the complex engineering involved in software systems for speech recognition and the effort that had been invested in building these systems on the basis of GMM-HMMs, the industry did not see a compelling argument for switching to neural networks. As a consequence, until the late 2000s, both academic and industrial research in using neural nets for speech recognition mostly focused on using neural nets to learn extra features for GMM-HMM systems.

Later, with **much larger and deeper models** and much larger datasets, recognition accuracy was dramatically improved by using neural networks to replace GMMs for the task of associating acoustic features to phonemes (or subphonemic states). Starting in 2009, speech researchers applied a form of deep learning based on unsupervised learning to speech recognition. This approach to deep learning was based on training undirected probabilistic models called restricted Boltzmann machines (RBMs) to model the input data. RBMs will be described in Part III. To solve speech recognition tasks, unsupervised pretraining was used to build deep feedforward networks whose layers were each initialized by training an RBM. These networks take spectral acoustic representations in a fixed-size input window (around a center frame) and predict the conditional probabilities of HMM states for that center frame. Training such deep networks helped to significantly improve the recognition rate on TIMIT (Mohamed *et al.*, 2009, 2012a), bringing down the phoneme error rate from about 26% to 20.7%. See Mohamed *et al.* (2012b) for an analysis of reasons for the success of these models. Extensions to the basic phone recognition pipeline included the addition of speaker-adaptive features (Mohamed *et al.*, 2011) that further reduced the error rate. This was quickly followed up by work to expand the architecture from phoneme recognition (which is what TIMIT is focused on) to large-vocabulary speech recognition (Dahl *et al.*, 2012), which involves not just recognizing phonemes but also recognizing sequences of words from a large vocabulary. Deep networks for speech recognition eventually shifted from being based on pretraining and Boltzmann machines to being based on techniques such as rectified linear units and dropout (Zeiler *et al.*, 2013; Dahl *et al.*, 2013). By that time, several of the major speech groups in industry had started exploring deep learning in collaboration with

academic researchers. Hinton *et al.* (2012a) describe the breakthroughs achieved by these collaborators, which are now deployed in products such as mobile phones.

Later, as these groups explored larger and larger labeled datasets and incorporated some of the methods for initializing, training, and setting up the architecture of deep nets, they realized that the unsupervised pretraining phase was either unnecessary or did not bring any significant improvement.

These breakthroughs in recognition performance for word error rate in speech recognition were unprecedented (around 30% improvement) and were following a long period of about ten years during which error rates did not improve much with the traditional GMM-HMM technology, in spite of the continuously growing size of training sets (see Fig. 2.4 of Deng and Yu (2014)). This created a rapid shift in the speech recognition community towards deep learning. In a matter of roughly two years, most of the industrial products for speech recognition incorporated deep neural networks and this success spurred a new wave of research into deep learning algorithms and architectures for ASR, which is still ongoing today.

One of these innovations was the use of convolutional networks (Sainath *et al.*, 2013) that replicate weights across time and frequency, improving over the earlier time-delay neural networks that replicated weights only across time. The new two-dimensional convolutional models regard the input spectrogram not as one long vector but as an image, with one axis corresponding to time and the other to frequency of spectral components.

Another important push, still ongoing, has been towards end-to-end deep learning speech recognition systems that completely remove the HMM. The first major breakthrough in this direction came from Graves *et al.* (2013) who trained a deep LSTM RNN (see Sec. 10.10), using MAP inference over the frame-to-phoneme alignment, as in LeCun *et al.* (1998b) and in the CTC framework (Graves *et al.*, 2006; Graves, 2012). A deep RNN (Graves *et al.*, 2013) has state variables from several layers at each time step, giving the unfolded graph two kinds of depth: ordinary depth due to a stack of layers, and depth due to time unfolding. This work brought the phoneme error rate on TIMIT to a record low of 17.7%. See Pascanu *et al.* (2014a) and Chung *et al.* (2014) for other variants of deep RNNs, applied in other settings.

Another contemporary step toward end-to-end deep learning ASR is to let the system learn how to “align” the acoustic-level information with the phonetic-level information (Chorowski *et al.*, 2014; Lu *et al.*, 2015).

## 12.4 Natural Language Processing

*Natural language processing* (NLP) is the use of human languages, such as English or French, by a computer. Computer programs typically read and emit specialized languages designed to allow efficient and unambiguous parsing by simple programs. More naturally occurring languages are often ambiguous and defy formal description. Natural language processing includes applications such as machine translation, in which the learner must read a sentence in one human language and emit an equivalent sentence in another human language. Many NLP applications are based on language models that define a probability distribution over sequences of words, characters or bytes in a natural language.

As with the other applications discussed in this chapter, very generic neural network techniques can be successfully applied to natural language processing. However, to achieve excellent performance and to scale well to large applications, some domain-specific strategies become important. To build an efficient model of natural language, we must usually use techniques that are specialized for processing sequential data. In many cases, we choose to regard natural language as a sequence of words, rather than a sequence of individual characters or bytes. Because the total number of possible words is so large, word-based language models must operate on an extremely high-dimensional and sparse discrete space. Several strategies have been developed to make models of such a space efficient, both in a computational and in a statistical sense.

### 12.4.1 *n*-grams

A *language model* defines a probability distribution over sequences of tokens in a natural language. Depending on how the model is designed, a token may be a word, a character, or even a byte. Tokens are always discrete entities. The earliest successful language models were based on models of fixed-length sequences of tokens called *n*-grams. An *n*-gram is a sequence of *n* tokens.

Models based on *n*-grams define the conditional probability of the *n*-th token given the preceding  $n - 1$  tokens. The model uses products of these conditional distributions to define the probability distribution over longer sequences:

$$P(x_1, \dots, x_\tau) = P(x_1, \dots, x_{n-1}) \prod_{t=n}^{\tau} P(x_t | x_{t-n+1}, \dots, x_{t-1}). \quad (12.5)$$

This decomposition is justified by the chain rule of probability. The probability distribution over the initial sequence  $P(x_1, \dots, x_{n-1})$  may be modeled by a different model with a smaller value of *n*.

Training  $n$ -gram models is straightforward because the maximum likelihood estimate can be computed simply by counting how many times each possible  $n$  gram occurs in the training set. Models based on  $n$ -grams have been the core building block of statistical language modeling for many decades ([Jelinek and Mercer, 1980](#); [Katz, 1987](#); [Chen and Goodman, 1999](#)).

For small values of  $n$ , models have particular names: *unigram* for  $n=1$ , *bigram* for  $n=2$ , and *trigram* for  $n=3$ . These names derive from the Latin prefixes for the corresponding numbers and the Greek suffix “-gram” denoting something that is written.

Usually we train both an  $n$ -gram model and an  $n-1$  gram model simultaneously. This makes it easy to compute

$$P(x_t | x_{t-n+1}, \dots, x_t) = \frac{P_n(x_{t-n+1}, \dots, x_t)}{P_{n-1}(x_{t-n+1}, \dots, x_{t-1})} \quad (12.6)$$

simply by looking up two stored probabilities. For this to exactly reproduce inference in  $P_n$ , we must omit the final character from each sequence when we train  $P_{n-1}$ .

As an example, we demonstrate how a trigram model computes the probability of the sentence “THE DOG RAN AWAY.” The first words of the sentence cannot be handled by the default formula based on conditional probability because there is no context at the beginning of the sentence. Instead, we must use the marginal probability over words at the start of the sentence. We thus evaluate  $P_3(\text{THE DOG RAN})$ . Finally, the last word may be predicted using the typical case, of using the conditional distribution  $P(\text{AWAY} | \text{DOG RAN})$ . Putting this together with Eq. 12.6, we obtain:

$$P(\text{THE DOG RAN AWAY}) = P_3(\text{THE DOG RAN})P_3(\text{DOG RAN AWAY})/P_2(\text{DOG RAN}). \quad (12.7)$$

A fundamental limitation of maximum likelihood for  $n$ -gram models is that  $P_n$  as estimated from training set counts is very likely to be zero in many cases, even though the tuple  $(x_{t-n+1}, \dots, x_t)$  may appear in the test set. This can cause two different kinds of catastrophic outcomes. When  $P_{n-1}$  is zero, the ratio is undefined, so the model does not even produce a sensible output. When  $P_{n-1}$  is non-zero but  $P_n$  is zero, the test log-likelihood is  $-\infty$ . To avoid such catastrophic outcomes, most  $n$ -gram models employ some form of *smoothing*. Smoothing techniques shift probability mass from the observed tuples to unobserved ones that are similar. See [Chen and Goodman \(1999\)](#) for a review and empirical comparisons. One basic technique consists of adding non-zero probability mass to all of the possible next symbol values. This method can be justified as Bayesian inference with a uniform

or Dirichlet prior over the count parameters. Another very popular idea is to form a mixture model containing higher-order and lower-order  $n$ -gram models, with the higher-order models providing more capacity and the lower-order models being more likely to avoid counts of zero. *Back-off methods* look-up the lower-order  $n$ -grams if the frequency of the context  $x_{t-1}, \dots, x_{t-n+1}$  is too small to use the higher-order model. More formally, they estimate the distribution over  $x_t$  by using contexts  $x_{t-n+k}, \dots, x_{t-1}$ , for increasing  $k$ , until a sufficiently reliable estimate is found.

Classical  $n$ -gram models are particularly vulnerable to the curse of dimensionality. There are  $|\mathbb{V}|^n$  possible  $n$ -grams and  $|\mathbb{V}|$  is often very large. Even with a massive training set and modest  $n$ , most  $n$ -grams will not occur in the training set. One way to view a classical  $n$ -gram model is that it is performing nearest-neighbor lookup. In other words, it can be viewed as a local non-parametric predictor, similar to  $k$ -nearest neighbors. The statistical problems facing these extremely local predictors are described in Sec. 5.11.2. The problem for a language model is even more severe than usual, because any two different words have the same distance from each other in one-hot vector space. It is thus difficult to leverage much information from any “neighbors”—only training examples that repeat literally the same context are useful for local generalization. To overcome these problems, a language model must be able to share knowledge between one word and other semantically similar words.

To improve the statistical efficiency of  $n$ -gram models, *class-based language models* (Brown *et al.*, 1992; Ney and Kneser, 1993; Niesler *et al.*, 1998) introduce the notion of word categories and then share statistical strength between words that are in the same category. The idea is to use a clustering algorithm to partition the set of words into clusters or classes, based on their co-occurrence frequencies with other words. The model can then use word class IDs rather than individual word IDs to represent the context on the right side of the conditioning bar. Composite models combining word-based and class-based models via mixing or back-off are also possible. Although word classes provide a way to generalize between sequences in which some word is replaced by another of the same class, much information is lost in this representation.

### 12.4.2 Neural Language Models

*Neural language models* or NLMs are a class of language model designed to overcome the curse of dimensionality problem for modeling natural language sequences by using a distributed representation of words (Bengio *et al.*, 2001). Unlike class-based  $n$ -gram models, neural language models are able to recognize that two words

are similar without losing the ability to encode each word as distinct from the other. Neural language models share statistical strength between one word (and its context) and other similar words and contexts. The distributed representation the model learns for each word enables this sharing by allowing the model to treat words that have features in common similarly. For example, if the word `dog` and the word `cat` map to representations that share many attributes, then sentences that contain the word `cat` can inform the predictions that will be made by the model for sentences that contain the word `dog`, and vice-versa. Because there are many such attributes, there are many ways in which generalization can happen, transferring information from each training sentence to an exponentially large number of semantically related sentences. The curse of dimensionality requires the model to generalize to a number of sentences that is exponential in the sentence length. The model counters this curse by relating each training sentence to an exponential number of similar sentences.

We sometimes call these word representations *word embeddings*. In this interpretation, we view the raw symbols as points in a space of dimension equal to the vocabulary size. The word representations embed those points in a feature space of lower dimension. In the original space, every word is represented by a one-hot vector, so every pair of words is at Euclidean distance  $\sqrt{2}$  from each other. In the embedding space, words that frequently appear in similar contexts (or any pair of words sharing some “features” learned by the model) are close to each other. This often results in words with similar meanings being neighbors. Fig. 12.3 zooms in on specific areas of a learned word embedding space to show how semantically similar words map to representations that are close to each other.

Neural networks in other domains also define embeddings. For example, a hidden layer of a convolutional network provides an “image embedding.” Usually NLP practitioners are much more interested in this idea of embeddings because natural language does not originally lie in a real-valued vector space. The hidden layer has provided a more qualitatively dramatic change in the way the data is represented.

The basic idea of using distributed representations to improve models for natural language processing is not restricted to neural networks. It may also be used with graphical models that have distributed representations in the form of multiple latent variables (Mnih and Hinton, 2007).

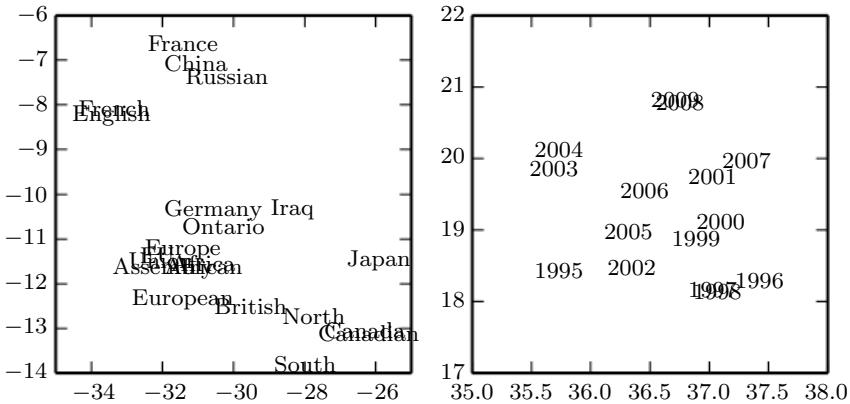


Figure 12.3: Two-dimensional visualizations of word embeddings obtained from a neural machine translation model (Bahdanau *et al.*, 2015), zooming in on specific areas where semantically related words have embedding vectors that are close to each other. Countries appear on the left and numbers on the right. Keep in mind that these embeddings are 2-D for the purpose of visualization. In real applications, embeddings typically have higher dimensionality and can simultaneously capture many kinds of similarity between words.

### 12.4.3 High-Dimensional Outputs

In many natural language applications, we often want our models to produce words (rather than characters) as the fundamental unit of the output. For large vocabularies, it can be very computationally expensive to represent an output distribution over the choice of a word, because the vocabulary size is large. In many applications,  $\mathbb{V}$  contains hundreds of thousands of words. The naive approach to representing such a distribution is to apply an affine transformation from a hidden representation to the output space, then apply the softmax function. Suppose we have a vocabulary  $\mathbb{V}$  with size  $|\mathbb{V}|$ . The weight matrix describing the linear component of this affine transformation is very large, because its output dimension is  $|\mathbb{V}|$ . This imposes a high memory cost to represent the matrix, and a high computational cost to multiply by it. Because the softmax is normalized across all  $|\mathbb{V}|$  outputs, it is necessary to perform the full matrix multiplication at training time as well as test time—we cannot calculate only the dot product with the weight vector for the correct output. The high computational costs of the output layer thus arise both at training time (to compute the likelihood and its gradient) and at test time (to compute probabilities for all or selected words). For specialized loss functions, the gradient can be computed efficiently (Vincent *et al.*, 2015), but the standard cross-entropy loss applied to a traditional softmax output layer poses

many difficulties.

Suppose that  $\mathbf{h}$  is the top hidden layer used to predict the output probabilities  $\hat{\mathbf{y}}$ . If we parametrize the transformation from  $\mathbf{h}$  to  $\hat{\mathbf{y}}$  with learned weights  $\mathbf{W}$  and learned biases  $\mathbf{b}$ , then the affine-softmax output layer performs the following computations:

$$a_i = b_i + \sum_j W_{ij} h_j \quad \forall i \in \{1, \dots, |\mathbb{V}|\}, \quad (12.8)$$

$$\hat{y}_i = \frac{e^{a_i}}{\sum_{i'=1}^{|\mathbb{V}|} e^{a_{i'}}}. \quad (12.9)$$

If  $\mathbf{h}$  contains  $n_h$  elements then the above operation is  $O(|\mathbb{V}|n_h)$ . With  $n_h$  in the thousands and  $|\mathbb{V}|$  in the hundreds of thousands, this operation dominates the computation of most neural language models.

#### 12.4.3.1 Use of a Short List

The first neural language models (Bengio *et al.*, 2001, 2003) dealt with the high cost of using a softmax over a large number of output words by limiting the vocabulary size to 10,000 or 20,000 words. Schwenk and Gauvain (2002) and Schwenk (2007) built upon this approach by splitting the vocabulary  $\mathbb{V}$  into a *shortlist*  $\mathbb{L}$  of most frequent words (handled by the neural net) and a tail  $\mathbb{T} = \mathbb{V} \setminus \mathbb{L}$  of more rare words (handled by an  $n$ -gram model). To be able to combine the two predictions, the neural net also has to predict the probability that a word appearing after context  $C$  belongs to the tail list. This may be achieved by adding an extra sigmoid output unit to provide an estimate of  $P(i \in \mathbb{T} | C)$ . The extra output can then be used to achieve an estimate of the probability distribution over all words in  $\mathbb{V}$  as follows:

$$P(y = i | C) = 1_{i \in \mathbb{L}} P(y = i | C, i \in \mathbb{L}) (1 - P(i \in \mathbb{T} | C)) \quad (12.10)$$

$$+ 1_{i \in \mathbb{T}} P(y = i | C, i \in \mathbb{T}) P(i \in \mathbb{T} | C) \quad (12.11)$$

where  $P(y = i | C, i \in \mathbb{L})$  is provided by the neural language model and  $P(y = i | C, i \in \mathbb{T})$  is provided by the  $n$ -gram model. With slight modification, this approach can also work using an extra output value in the neural language model's softmax layer, rather than a separate sigmoid unit.

An obvious disadvantage of the short list approach is that the potential generalization advantage of the neural language models is limited to the most frequent words, where, arguably, it is the least useful. This disadvantage has stimulated the exploration of alternative methods to deal with high-dimensional outputs, described below.

### 12.4.3.2 Hierarchical Softmax

A classical approach (Goodman, 2001) to reducing the computational burden of high-dimensional output layers over large vocabulary sets  $\mathbb{V}$  is to decompose probabilities hierarchically. Instead of necessitating a number of computations proportional to  $|\mathbb{V}|$  (and also proportional to the number of hidden units,  $n_h$ ), the  $|\mathbb{V}|$  factor can be reduced to as low as  $\log |\mathbb{V}|$ . Bengio (2002) and Morin and Bengio (2005) introduced this factorized approach to the context of neural language models.

One can think of this hierarchy as building categories of words, then categories of categories of words, then categories of categories of categories of words, etc. These nested categories form a tree, with words at the leaves. In a balanced tree, the tree has depth  $O(\log |\mathbb{V}|)$ . The probability of choosing a word is given by the product of the probabilities of choosing the branch leading to that word at every node on a path from the root of the tree to the leaf containing the word. Fig. 12.4 illustrates a simple example. Mnih and Hinton (2009) also describe how to use multiple paths to identify a single word in order to better model words that have multiple meanings. Computing the probability of a word then involves summation over all of the paths that lead to that word.

To predict the conditional probabilities required at each node of the tree, we typically use a logistic regression model at each node of the tree, and provide the same context  $C$  as input to all of these models. Because the correct output is encoded in the training set, we can use supervised learning to train the logistic regression models. This is typically done using a standard cross-entropy loss, corresponding to maximizing the log-likelihood of the correct sequence of decisions.

Because the output log-likelihood can be computed efficiently (as low as  $\log |\mathbb{V}|$  rather than  $|\mathbb{V}|$ ), its gradients may also be computed efficiently. This includes not only the gradient with respect to the output parameters but also the gradients with respect to the hidden layer activations.

It is possible but usually not practical to optimize the tree structure to minimize the expected number of computations. Tools from information theory specify how to choose the optimal binary code given the relative frequencies of the words. To do so, we could structure the tree so that the number of bits associated with a word is approximately equal to the logarithm of the frequency of that word. However, in practice, the computational savings are typically not worth the effort because the computation of the output probabilities is only one part of the total computation in the neural language model. For example, suppose there are  $l$  fully connected hidden layers of width  $n_h$ . Let  $n_b$  be the weighted average of the number of bits

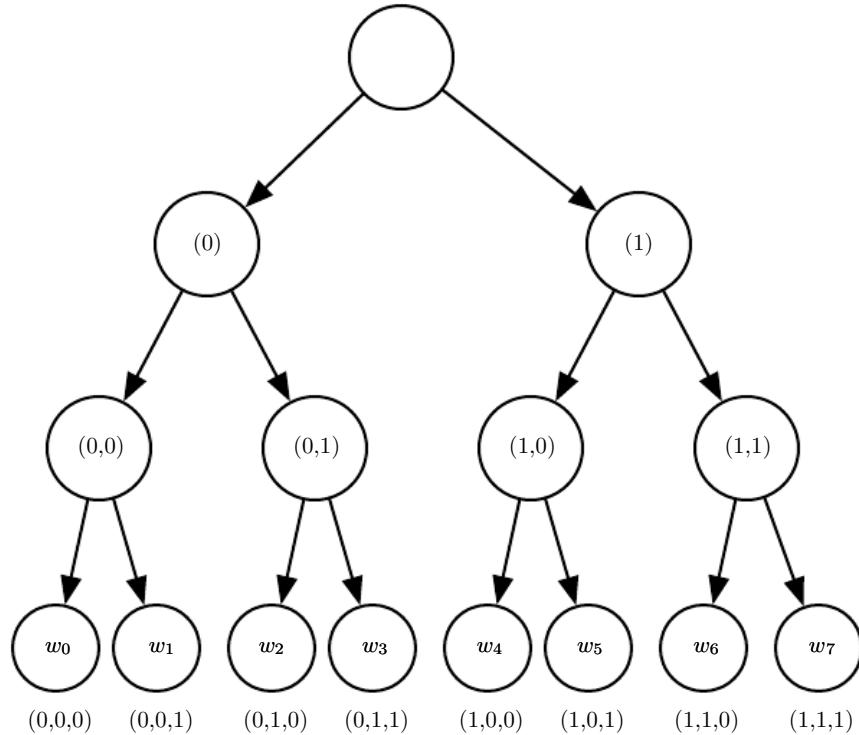


Figure 12.4: Illustration of a simple hierarchy of word categories, with 8 words  $w_0, \dots, w_7$  organized into a three level hierarchy. The leaves of the tree represent actual specific words. Internal nodes represent groups of words. Any node can be indexed by the sequence of binary decisions (0=left, 1=right) to reach the node from the root. Super-class (0) contains the classes (0,0) and (0,1), which respectively contain the sets of words  $\{w_0, w_1\}$  and  $\{w_2, w_3\}$ , and similarly super-class (1) contains the classes (1,0) and (1,1), which respectively contain the words  $(w_4, w_5)$  and  $(w_6, w_7)$ . If the tree is sufficiently balanced, the maximum depth (number of binary decisions) is on the order of the logarithm of the number of words  $|\mathbb{V}|$ : the choice of one out of  $|\mathbb{V}|$  words can be obtained by doing  $O(\log |\mathbb{V}|)$  operations (one for each of the nodes on the path from the root). In this example, computing the probability of a word  $y$  can be done by multiplying three probabilities, associated with the binary decisions to move left or right at each node on the path from the root to a node  $y$ . Let  $b_i(y)$  be the  $i$ -th binary decision when traversing the tree towards the value  $y$ . The probability of sampling an output  $y$  decomposes into a product of conditional probabilities, using the chain rule for conditional probabilities, with each node indexed by the prefix of these bits. For example, node (1,0) corresponds to the prefix  $(b_0(w_4) = 1, b_1(w_4) = 0)$ , and the probability of  $w_4$  can be decomposed as follows:

$$P(y = w_4) = P(b_0 = 1, b_1 = 0, b_2 = 0) \quad (12.12)$$

$$= P(b_0 = 1)P(b_1 = 0 | b_0 = 1)P(b_2 = 0 | b_0 = 1, b_1 = 0). \quad (12.13)$$

required to identify a word, with the weighting given by the frequency of these words. In this example, the number of operations needed to compute the hidden activations grows as  $O(\ln_h^2)$  while the output computations grow as  $O(n_h n_b)$ . As long as  $n_b \leq \ln_h$ , we can reduce computation more by shrinking  $n_h$  than by shrinking  $n_b$ . Indeed,  $n_b$  is often small. Because the size of the vocabulary rarely exceeds a million words and  $\log_2(10^6) \approx 20$ , it is possible to reduce  $n_b$  to about 20, but  $n_h$  is often much larger, around  $10^3$  or more. Rather than carefully optimizing a tree with a branching factor of 2, one can instead define a tree with depth two and a branching factor of  $\sqrt{|\mathbb{V}|}$ . Such a tree corresponds to simply defining a set of mutually exclusive word classes. The simple approach based on a tree of depth two captures most of the computational benefit of the hierarchical strategy.

One question that remains somewhat open is how to best define these word classes, or how to define the word hierarchy in general. Early work used existing hierarchies (Morin and Bengio, 2005) but the hierarchy can also be learned, ideally jointly with the neural language model. Learning the hierarchy is difficult. An exact optimization of the log-likelihood appears intractable because the choice of a word hierarchy is a discrete one, not amenable to gradient-based optimization. However, one could use discrete optimization to approximately optimize the partition of words into word classes.

An important advantage of the hierarchical softmax is that it brings computational benefits both at training time and at test time, if at test time we want to compute the probability of specific words.

Of course, computing the probability of all  $|\mathbb{V}|$  words will remain expensive even with the hierarchical softmax. Another important operation is selecting the most likely word in a given context. Unfortunately the tree structure does not provide an efficient and exact solution to this problem.

A disadvantage is that in practice the hierarchical softmax tends to give worse test results than sampling-based methods we will describe next. This may be due to a poor choice of word classes.

#### 12.4.3.3 Importance Sampling

One way to speed up the training of neural language models is to avoid explicitly computing the contribution of the gradient from all of the words that do not appear in the next position. Every incorrect word should have low probability under the model. It can be computationally costly to enumerate all of these words. Instead, it is possible to sample only a subset of the words. Using the notation introduced

in Eq. 12.8, the gradient can be written as follows:

$$\frac{\partial \log P(y | C)}{\partial \theta} = \frac{\partial \log \text{softmax}_y(\mathbf{a})}{\partial \theta} \quad (12.14)$$

$$= \frac{\partial}{\partial \theta} \log \frac{e^{a_y}}{\sum_i e^{a_i}} \quad (12.15)$$

$$= \frac{\partial}{\partial \theta} (a_y - \log \sum_i e^{a_i}) \quad (12.16)$$

$$= \frac{\partial a_y}{\partial \theta} - \sum_i P(y = i | C) \frac{\partial a_i}{\partial \theta} \quad (12.17)$$

where  $\mathbf{a}$  is the vector of pre-softmax activations (or scores), with one element per word. The first term is the *positive phase* term (pushing  $a_y$  up) while the second term is the *negative phase* term (pushing  $a_i$  down for all  $i$ , with weight  $P(i | C)$ ). Since the negative phase term is an expectation, we can estimate it with a Monte Carlo sample. However, that would require sampling from the model itself. Sampling from the model requires computing  $P(i | C)$  for all  $i$  in the vocabulary, which is precisely what we are trying to avoid.

Instead of sampling from the model, one can sample from another distribution, called the proposal distribution (denoted  $q$ ), and use appropriate weights to correct for the bias introduced by sampling from the wrong distribution (Bengio and Sénecal, 2003; Bengio and Sénecal, 2008). This is an application of a more general technique called *importance sampling*, which will be described in more detail in Sec. 17.2. Unfortunately, even exact importance sampling is not efficient because it requires computing weights  $p_i/q_i$ , where  $p_i = P(i | C)$ , which can only be computed if all the scores  $a_i$  are computed. The solution adopted for this application is called *biased importance sampling*, where the importance weights are normalized to sum to 1. When negative word  $n_i$  is sampled, the associated gradient is weighted by

$$w_i = \frac{p_{n_i}/q_{n_i}}{\sum_{j=1}^N p_{n_j}/q_{n_j}}. \quad (12.18)$$

These weights are used to give the appropriate importance to the  $m$  negative samples from  $q$  used to form the estimated negative phase contribution to the gradient:

$$\sum_{i=1}^{|V|} P(i | C) \frac{\partial a_i}{\partial \theta} \approx \frac{1}{m} \sum_{i=1}^m w_i \frac{\partial a_{n_i}}{\partial \theta}. \quad (12.19)$$

A unigram or a bigram distribution works well as the proposal distribution  $q$ . It is easy to estimate the parameters of such a distribution from data. After estimating the parameters, it is also possible to sample from such a distribution very efficiently.

Importance sampling is not only useful for speeding up models with large softmax outputs. More generally, it is useful for accelerating training with large sparse output layers, where the output is a sparse vector rather than a 1-of- $n$  choice. An example is a *bag of words*. A bag of words is a sparse vector  $\mathbf{v}$  where  $v_i$  indicates the presence or absence of word  $i$  from the vocabulary in the document. Alternately,  $v_i$  can indicate the number of times that word  $i$  appears. Machine learning models that emit such sparse vectors can be expensive to train for a variety of reasons. Early in learning, the model may not actually choose to make the output truly sparse. Moreover, the loss function we use for training might most naturally be described in terms of comparing every element of the output to every element of the target. This means that it is not always clear that there is a computational benefit to using sparse outputs, because the model may choose to make the majority of the output non-zero and all of these non-zero values need to be compared to the corresponding training target, even if the training target is zero. [Dauphin et al. \(2011\)](#) demonstrated that such models can be accelerated using importance sampling. The efficient algorithm minimizes the loss reconstruction for the “positive words” (those that are non-zero in the target) and an equal number of “negative words.” The negative words are chosen randomly, using a heuristic to sample words that are more likely to be mistaken. The bias introduced by this heuristic oversampling can then be corrected using importance weights.

In all of these cases, the computational complexity of gradient estimation for the output layer is reduced to be proportional to the number of negative samples rather than proportional to the size of the output vector.

#### 12.4.3.4 Noise-Contrastive Estimation and Ranking Loss

Other approaches based on sampling have been proposed to reduce the computational cost of training neural language models with large vocabularies. An early example is the ranking loss proposed by [Collobert and Weston \(2008a\)](#), which views the output of the neural language model for each word as a score and tries to make the score of the correct word  $a_y$  be ranked high in comparison to the other scores  $a_i$ . The ranking loss proposed then is

$$L = \sum_i \max(0, 1 - a_y + a_i). \quad (12.20)$$

The gradient is zero for the  $i$ -th term if the score of the observed word,  $a_y$ , is greater than the score of the negative word  $a_i$  by a margin of 1. One issue with this criterion is that it does not provide estimated conditional probabilities, which

are useful in some applications, including speech recognition and text generation (including conditional text generation tasks such as translation).

A more recently used training objective for neural language model is noise-contrastive estimation, which is introduced in Sec. 18.6. This approach has been successfully applied to neural language models (Mnih and Teh, 2012; Mnih and Kavukcuoglu, 2013).

#### 12.4.4 Combining Neural Language Models with $n$ -grams

A major advantage of  $n$ -gram models over neural networks is that  $n$ -gram models achieve high model capacity (by storing the frequencies of very many tuples) while requiring very little computation to process an example (by looking up only a few tuples that match the current context). If we use hash tables or trees to access the counts, the computation used for  $n$ -grams is almost independent of capacity. In comparison, doubling a neural network’s number of parameters typically also roughly doubles its computation time. Exceptions include models that avoid using all parameters on each pass. Embedding layers index only a single embedding in each pass, so we can increase the vocabulary size without increasing the computation time per example. Some other models, such as tiled convolutional networks, can add parameters while reducing the degree of parameter sharing in order to maintain the same amount of computation. However, typical neural network layers based on matrix multiplication use an amount of computation proportional to the number of parameters.

One easy way to add capacity is thus to combine both approaches in an ensemble consisting of a neural language model and an  $n$ -gram language model (Bengio *et al.*, 2001, 2003). As with any ensemble, this technique can reduce test error if the ensemble members make independent mistakes. The field of ensemble learning provides many ways of combining the ensemble members’ predictions, including uniform weighting and weights chosen on a validation set. Mikolov *et al.* (2011a) extended the ensemble to include not just two models but a large array of models. It is also possible to pair a neural network with a maximum entropy model and train both jointly (Mikolov *et al.*, 2011b). This approach can be viewed as training a neural network with an extra set of inputs that are connected directly to the output, and not connected to any other part of the model. The extra inputs are indicators for the presence of particular  $n$ -grams in the input context, so these variables are very high-dimensional and very sparse. The increase in model capacity is huge—the new portion of the architecture contains up to  $|sV|^n$  parameters—but the amount of added computation needed to process an input is minimal because the extra inputs are very sparse.

### 12.4.5 Neural Machine Translation

Machine translation is the task of reading a sentence in one natural language and emitting a sentence with the equivalent meaning in another language. Machine translation systems often involve many components. At a high level, there is often one component that proposes many candidate translations. Many of these translations will not be grammatical due to differences between the languages. For example, many languages put adjectives after nouns, so when translated to English directly they yield phrases such as “apple red.” The proposal mechanism suggests many variants of the suggested translation, ideally including “red apple.” A second component of the translation system, a language model, evaluates the proposed translations, and can score “red apple” as better than “apple red.”

The earliest use of neural networks for machine translation was to upgrade the language model of a translation system by using a neural language model ([Schwenk et al., 2006](#); [Schwenk, 2010](#)). Previously, most machine translation systems had used an  $n$ -gram model for this component. The  $n$ -gram based models used for machine translation include not just traditional back-off  $n$ -gram models ([Jelinek and Mercer, 1980](#); [Katz, 1987](#); [Chen and Goodman, 1999](#)) but also *maximum entropy language models* ([Berger et al., 1996](#)), in which an affine-softmax layer predicts the next word given the presence of frequent  $n$ -grams in the context.

Traditional language models simply report the probability of a natural language sentence. Because machine translation involves producing an output sentence given an input sentence, it makes sense to extend the natural language model to be conditional. As described in Sec. [6.2.1.1](#), it is straightforward to extend a model that defines a marginal distribution over some variable to define a conditional distribution over that variable given a context  $C$ , where  $C$  might be a single variable or a list of variables. [Devlin et al. \(2014\)](#) beat the state-of-the-art in some statistical machine translation benchmarks by using an MLP to score a phrase  $t_1, t_2, \dots, t_k$  in the target language given a phrase  $s_1, s_2, \dots, s_n$  in the source language. The MLP estimates  $P(t_1, t_2, \dots, t_k | s_1, s_2, \dots, s_n)$ . The estimate formed by this MLP replaces the estimate provided by conditional  $n$ -gram models.

A drawback of the MLP-based approach is that it requires the sequences to be preprocessed to be of fixed length. To make the translation more flexible, we would like to use a model that can accommodate variable length inputs and variable length outputs. An RNN provides this ability. Sec. [10.2.4](#) describes several ways of constructing an RNN that represents a conditional distribution over a sequence given some input, and Sec. [10.4](#) describes how to accomplish this conditioning when the input is a sequence. In all cases, one model first reads the input sequence and emits a data structure that summarizes the input sequence. We call this

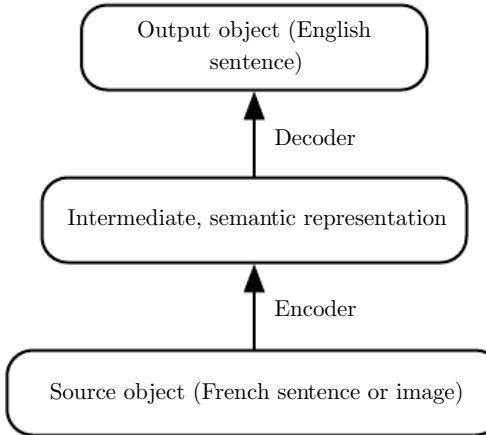


Figure 12.5: The encoder-decoder architecture to map back and forth between a surface representation (such as a sequence of words or an image) and a semantic representation. By using the output of an encoder of data from one modality (such as the encoder mapping from French sentences to hidden representations capturing the meaning of sentences) as the input to a decoder for another modality (such as the decoder mapping from hidden representations capturing the meaning of sentences to English), we can train systems that translate from one modality to another. This idea has been applied successfully not just to machine translation but also to caption generation from images.

summary the “context”  $C$ . The context  $C$  may be a list of vectors, or it may be a vector or tensor. The model that reads the input to produce  $C$  may be an RNN (Cho *et al.*, 2014a; Sutskever *et al.*, 2014; Jean *et al.*, 2014) or a convolutional network (Kalchbrenner and Blunsom, 2013). A second model, usually an RNN, then reads the context  $C$  and generates a sentence in the target language. This general idea of an encoder-decoder framework for machine translation is illustrated in Fig. 12.5.

In order to generate an entire sentence conditioned on the source sentence, the model must have a way to represent the entire source sentence. Earlier models were only able to represent individual words or phrases. From a representation learning point of view, it can be useful to learn a representation in which sentences that have the same meaning have similar representations regardless of whether they were written in the source language or the target language. This strategy was explored first using a combination of convolutions and RNNs (Kalchbrenner and Blunsom, 2013). Later work introduced the use of an RNN for scoring proposed translations (Cho *et al.*, 2014a) and for generating translated sentences (Sutskever *et al.*, 2014). Jean *et al.* (2014) scaled these models to larger vocabularies.

### 12.4.5.1 Using an Attention Mechanism and Aligning Pieces of Data

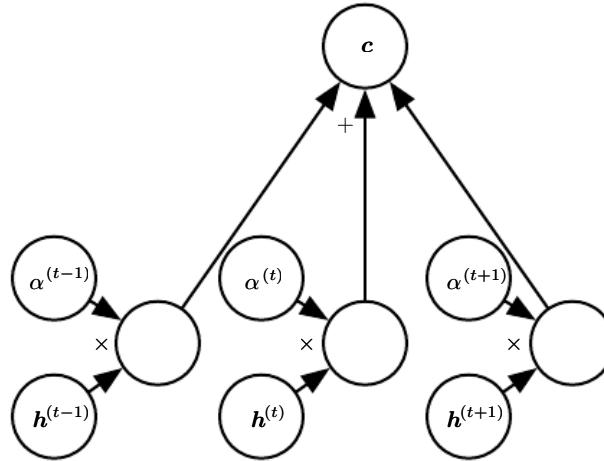


Figure 12.6: A modern attention mechanism, as introduced by [Bahdanau et al. \(2015\)](#), is essentially a weighted average. A context vector  $c$  is formed by taking a weighted average of feature vectors  $\mathbf{h}^{(t)}$  with weights  $\alpha^{(t)}$ . In some applications, the feature vectors  $\mathbf{h}$  are hidden units of a neural network, but they may also be raw input to the model. The weights  $\alpha^{(t)}$  are produced by the model itself. They are usually values in the interval  $[0, 1]$  and are intended to concentrate around just one  $\mathbf{h}^{(t)}$  so that the weighted average approximates reading that one specific time step precisely. The weights  $\alpha^{(t)}$  are usually produced by applying a softmax function to relevance scores emitted by another portion of the model. The attention mechanism is more expensive computationally than directly indexing the desired  $\mathbf{h}^{(t)}$ , but direct indexing cannot be trained with gradient descent. The attention mechanism based on weighted averages is a smooth, differentiable approximation that can be trained with existing optimization algorithms.

Using a fixed-size representation to capture all the semantic details of a very long sentence of say 60 words is very difficult. It can be achieved by training a sufficiently large RNN well enough and for long enough, as demonstrated by [Cho et al. \(2014a\)](#) and [Sutskever et al. \(2014\)](#). However, a more efficient approach is to read the whole sentence or paragraph (to get the context and the gist of what is being expressed), then produce the translated words one at a time, each time focusing on a different part of the input sentence in order to gather the semantic details that are required to produce the next output word. That is exactly the idea that [Bahdanau et al. \(2015\)](#) first introduced. The attention mechanism used to focus on specific parts of the input sequence at each time step is illustrated in Fig. 12.6.

We can think of an attention-based system as having three components:

1. A process that “*reads*” raw data (such as source words in a source sentence), and converts them into distributed representations, with one feature vector associated with each word position.
2. A list of feature vectors storing the output of the reader. This can be understood as a “*memory*” containing a sequence of facts, which can be retrieved later, not necessarily in the same order, without having to visit all of them.
3. A process that “*exploits*” the content of the memory to sequentially perform a task, at each time step having the ability put attention on the content of one memory element (or a few, with a different weight).

The third component generates the translated sentence.

When words in a sentence written in one language are aligned with corresponding words in a translated sentence in another language, it becomes possible to relate the corresponding word embeddings. Earlier work showed that one could learn a kind of translation matrix relating the word embeddings in one language with the word embeddings in another (Kočiský *et al.*, 2014), yielding lower alignment error rates than traditional approaches based on the frequency counts in the phrase table. There is even earlier work on learning cross-lingual word vectors (Klementiev *et al.*, 2012). Many extensions to this approach are possible. For example, more efficient cross-lingual alignment (Gouws *et al.*, 2014) allows training on larger datasets.

#### 12.4.6 Historical Perspective

The idea of distributed representations for symbols was introduced by Rumelhart *et al.* (1986a) in one of the first explorations of back-propagation, with symbols corresponding to the identity of family members and the neural network capturing the relationships between family members, with training examples forming triplets such as (Colin, Mother, Victoria). The first layer of the neural network learned a representation of each family member. For example, the features for Colin might represent which family tree Colin was in, what branch of that tree he was in, what generation he was from, etc. One can think of the neural network as computing learned rules relating these attributes together in order to obtain the desired predictions. The model can then make predictions such as inferring who is the mother of Colin.

The idea of forming an embedding for a symbol was extended to the idea of an embedding for a word by Deerwester *et al.* (1990). These embeddings were learned using the SVD. Later, embeddings would be learned by neural networks.

The history of natural language processing is marked by transitions in the popularity of different ways of representing the input to the model. Following this early work on symbols or words, some of the earliest applications of neural networks to NLP (Mikkelainen and Dyer, 1991; Schmidhuber, 1996) represented the input as a sequence of characters.

Bengio *et al.* (2001) returned the focus to modeling words and introduced neural language models, which produce interpretable word embeddings. These neural models have scaled up from defining representations of a small set of symbols in the 1980s to millions of words (including proper nouns and misspellings) in modern applications. This computational scaling effort led to the invention of the techniques described above in Sec. 12.4.3.

Initially, the use of words as the fundamental units of language models yielded improved language modeling performance (Bengio *et al.*, 2001). To this day, new techniques continually push both character-based models (Sutskever *et al.*, 2011) and word-based models forward, with recent work (Gillick *et al.*, 2015) even modeling individual bytes of Unicode characters.

The ideas behind neural language models have been extended into several natural language processing applications, such as parsing (Henderson, 2003, 2004; Collobert, 2011), part-of-speech tagging, semantic role labeling, chunking, etc, sometimes using a single multi-task learning architecture (Collobert and Weston, 2008a; Collobert *et al.*, 2011a) in which the word embeddings are shared across tasks.

Two-dimensional visualizations of embeddings became a popular tool for analyzing language models following the development of the t-SNE dimensionality reduction algorithm (van der Maaten and Hinton, 2008) and its high-profile application to visualization word embeddings by Joseph Turian in 2009.

## 12.5 Other Applications

In this section we cover a few other types of applications of deep learning that are different from the standard object recognition, speech recognition and natural language processing tasks discussed above. Part III of this book will expand that scope even further to include tasks requiring the ability to generate rich high-dimensional samples (unlike “the next word,” in language models).

### 12.5.1 Recommender Systems

One of the major families of applications of machine learning in the information technology sector is the ability to make recommendations of items to potential users or customers. Two major types of applications can be distinguished: online advertising and item recommendations (often these recommendations are still for the purpose of selling a product). Both rely on predicting the association between a user and an item, either to predict the probability of some action (the user buying the product, or some proxy for this action) or the expected gain (which may depend on the value of the product) if an ad is shown or a recommendation is made regarding that product to that user. The internet is currently financed in great part by various forms of online advertising. There are major parts of the economy that rely on online shopping. Companies including Amazon and eBay use machine learning, including deep learning, for their product recommendations. Sometimes, the items are not products that are actually for sale. Examples include selecting posts to display on social network news feeds, recommending movies to watch, recommending jokes, recommending advice from experts, matching players for video games, or matching people in dating services.

Often, this association problem is handled like a supervised learning problem: given some information about the item and about the user, predict the proxy of interest (user clicks on ad, user enters a rating, user clicks on a “like” button, user buys product, user spends some amount of money on the product, user spends time visiting a page for the product, etc). This often ends up being either a regression problem (predicting some conditional expected value) or a probabilistic classification problem (predicting the conditional probability of some discrete event).

The early work on recommender systems relied on minimal information as inputs for these predictions: the user ID and the item ID. In this context, the only way to generalize is to rely on the similarity between the patterns of values of the target variable for different users or for different items. Suppose that user 1 and user 2 both like items A, B and C. From this, we may infer that user 1 and user 2 have similar tastes. If user 1 likes item D, then this should be a strong cue that user 2 will also like D. Algorithms based on this principle come under the name of *collaborative filtering*. Both non-parametric approaches (such as nearest-neighbor methods based on the estimated similarity between patterns of preferences) and parametric methods are possible. Parametric methods often rely on learning a distributed representation (also called an embedding) for each user and for each item. Bilinear prediction of the target variable (such as a rating) is a simple parametric method that is highly successful and often found as a component of

state-of-the-art systems. The prediction is obtained by the dot product between the user embedding and the item embedding (possibly corrected by constants that depend only on either the user ID or the item ID). Let  $\hat{\mathbf{R}}$  be the matrix containing our predictions,  $\mathbf{A}$  a matrix with user embeddings in its rows and  $\mathbf{B}$  a matrix with item embeddings in its columns. Let  $\mathbf{b}$  and  $\mathbf{c}$  be vectors that contain respectively a kind of bias for each user (representing how grumpy or positive that user is in general) and for each item (representing its general popularity). The bilinear prediction is thus obtained as follows:

$$\hat{R}_{u,i} = b_u + c_i + \sum_j A_{u,j} B_{j,i}. \quad (12.21)$$

Typically one wants to minimize the squared error between predicted ratings  $\hat{R}_{u,i}$  and actual ratings  $R_{u,i}$ . User embeddings and item embeddings can then be conveniently visualized when they are first reduced to a low dimension (two or three), or they can be used to compare users or items against each other, just like word embeddings. One way to obtain these embeddings is by performing a singular value decomposition of the matrix  $\mathbf{R}$  of actual targets (such as ratings). This corresponds to factorizing  $\mathbf{R} = \mathbf{UDV}'$  (or a normalized variant) into the product of two factors, the lower rank matrices  $\mathbf{A} = \mathbf{UD}$  and  $\mathbf{B} = \mathbf{V}'$ . One problem with the SVD is that it treats the missing entries in an arbitrary way, as if they corresponded to a target value of 0. Instead we would like to avoid paying any cost for the predictions made on missing entries. Fortunately, the sum of squared errors on the observed ratings can also be easily minimized by gradient-based optimization. The SVD and the bilinear prediction of Eq. 12.21 both performed very well in the competition for the Netflix prize (Bennett and Lanning, 2007), aiming at predicting ratings for films, based only on previous ratings by a large set of anonymous users. Many machine learning experts participated in this competition, which took place between 2006 and 2009. It raised the level of research in recommender systems using advanced machine learning and yielded improvements in recommender systems. Even though it did not win by itself, the simple bilinear prediction or SVD was a component of the ensemble models presented by most of the competitors, including the winners (Töscher *et al.*, 2009; Koren, 2009).

Beyond these bilinear models with distributed representations, one of the first uses of neural networks for collaborative filtering is based on the RBM undirected probabilistic model (Salakhutdinov *et al.*, 2007). RBMs were an important element of the ensemble of methods that won the Netflix competition (Töscher *et al.*, 2009; Koren, 2009). More advanced variants on the idea of factorizing the ratings matrix have also been explored in the neural networks community (Salakhutdinov and

Mnih, 2008).

However, there is a basic limitation of collaborative filtering systems: when a new item or a new user is introduced, its lack of rating history means that there is no way to evaluate its similarity with other items or users (respectively), or the degree of association between, say, that new user and existing items. This is called the problem of cold-start recommendations. A general way of solving the cold-start recommendation problem is to introduce extra information about the individual users and items. For example, this extra information could be user profile information or features of each item. Systems that use such information are called *content-based recommender systems*. The mapping from a rich set of user features or item features to an embedding can be learned through a deep learning architecture (Huang *et al.*, 2013; Elkahky *et al.*, 2015).

Specialized deep learning architectures such as convolutional networks have also been applied to learn to extract features from rich content such as from musical audio tracks, for music recommendation (van den Oörd *et al.*, 2013). In that work, the convolutional net takes acoustic features as input and computes an embedding for the associated song. The dot product between this song embedding and the embedding for a user is then used to predict whether a user will listen to the song.

### 12.5.1.1 Exploration Versus Exploitation

When making recommendations to users, an issue arises that goes beyond ordinary supervised learning and into the realm of reinforcement learning. Many recommendation problems are most accurately described theoretically as *contextual bandits* (Langford and Zhang, 2008; Lu *et al.*, 2010). The issue is that when we use the recommendation system to collect data, we get a biased and incomplete view of the preferences of users: we only see the responses of users to the items they were recommended and not to the other items. In addition, in some cases we may not get any information on users for whom no recommendation has been made (for example, with ad auctions, it may be that the price proposed for an ad was below a minimum price threshold, or does not win the auction, so the ad is not shown at all). More importantly, we get no information about what outcome would have resulted from recommending any of the other items. This would be like training a classifier by picking one class  $\hat{y}$  for each training example  $\mathbf{x}$  (typically the class with the highest probability according to the model) and then only getting as feedback whether this was the correct class or not. Clearly, each example conveys less information than in the supervised case where the true label  $y$  is directly accessible, so more examples are necessary. Worse, if we are not careful, we could end up with a system that continues picking the wrong decisions even as more

and more data is collected, because the correct decision initially had a very low probability: until the learner picks that correct decision, it does not learn about the correct decision. This is similar to the situation in reinforcement learning where only the reward for the selected action is observed. In general, reinforcement learning can involve a sequence of many actions and many rewards. The bandits scenario is a special case of reinforcement learning, in which the learner takes only a single action and receives a single reward. The bandit problem is easier in the sense that the learner knows which reward is associated with which action. In the general reinforcement learning scenario, a high reward or a low reward might have been caused by a recent action or by an action in the distant past. The term *contextual* bandits refers to the case where the action is taken in the context of some input variable that can inform the decision. For example, we at least know the user identity, and we want to pick an item. The mapping from context to action is also called a *policy*. The feedback loop between the learner and the data distribution (which now depends on the actions of the learner) is a central research issue in the reinforcement learning and bandits literature.

Reinforcement learning requires choosing a tradeoff between *exploration* and *exploitation*. Exploitation refers to taking actions that come from the current, best version of the learned policy—actions that we know will achieve a high reward. Exploration refers to taking actions specifically in order to obtain more training data. If we know that given context  $x$ , action  $a$  gives us a reward of 1, we do not know whether that is the best possible reward. We may want to exploit our current policy and continue taking action  $a$  in order to be relatively sure of obtaining a reward of 1. However, we may also want to explore by trying action  $a'$ . We do not know what will happen if we try action  $a'$ . We hope to get a reward of 2, but we run the risk of getting a reward of 0. Either way, we at least gain some knowledge.

Exploration can be implemented in many ways, ranging from occasionally taking random actions intended to cover the entire space of possible actions, to model-based approaches that compute a choice of action based on its expected reward and the model’s amount of uncertainty about that reward.

Many factors determine the extent to which we prefer exploration or exploitation. One of the most prominent factors is the time scale we are interested in. If the agent has only a short amount of time to accrue reward, then we prefer more exploitation. If the agent has a long time to accrue reward, then we begin with more exploration so that future actions can be planned more effectively with more knowledge. As time progresses and our learned policy improves, we move toward more exploitation.

Supervised learning has no tradeoff between exploration and exploitation

because the supervision signal always specifies which output is correct for each input. There is no need to try out different outputs to determine if one is better than the model’s current output—we always know that the label is the best output.

Another difficulty arising in the context of reinforcement learning, besides the exploration-exploitation trade-off, is the difficulty of evaluating and comparing different policies. Reinforcement learning involves interaction between the learner and the environment. This feedback loop means that it is not straightforward to evaluate the learner’s performance using a fixed set of test set input values. The policy itself determines which inputs will be seen. [Dudik et al. \(2011\)](#) present techniques for evaluating contextual bandits.

## 12.5.2 Knowledge Representation, Reasoning and Question Answering

Deep learning approaches have been very successful in language modeling, machine translation and natural language processing due to the use of embeddings for symbols ([Rumelhart et al., 1986a](#)) and words ([Deerwester et al., 1990; Bengio et al., 2001](#)). These embeddings represent semantic knowledge about individual words and concepts. A research frontier is to develop embeddings for phrases and for relations between words and facts. Search engines already use machine learning for this purpose but much more remains to be done to improve these more advanced representations.

### 12.5.2.1 Knowledge, Relations and Question Answering

indexRelations One interesting research direction is determining how distributed representations can be trained to capture the *relations* between two entities. These relations allow us to formalize facts about objects and how objects interact with each other.

In mathematics, a *binary relation* is a set of ordered pairs of objects. Pairs that are in the set are said to have the relation while those who are not in the set do not. For example, we can define the relation “is less than” on the set of entities  $\{1, 2, 3\}$  by defining the set of ordered pairs  $\mathbb{S} = \{(1, 2), (1, 3), (2, 3)\}$ . Once this relation is defined, we can use it like a verb. Because  $(1, 2) \in \mathbb{S}$ , we say that 1 is less than 2. Because  $(2, 1) \notin \mathbb{S}$ , we can not say that 2 is less than 1. Of course, the entities that are related to one another need not be numbers. We could define a relation `is_a_type_of` containing tuples like `(dog, mammal)`.

In the context of AI, we think of a relation as a sentence in a syntactically

simple and highly structured language. The relation plays the role of a verb, while two arguments to the relation play the role of its subject and object. These sentences take the form of a triplet of tokens

$$(\text{subject}, \text{verb}, \text{object}) \quad (12.22)$$

with values

$$(\text{entity}_i, \text{relation}_j, \text{entity}_k). \quad (12.23)$$

We can also define an *attribute*, a concept analogous to a relation, but taking only one argument:

$$(\text{entity}_i, \text{attribute}_j). \quad (12.24)$$

For example, we could define the `has_fur` attribute, and apply it to entities like `dog`.

Many applications require representing relations and reasoning about them. How should we best do this within the context of neural networks?

Machine learning models of course require training data. We can infer relations between entities from training datasets consisting of unstructured natural language. There are also structured databases that identify relations explicitly. A common structure for these databases is the *relational database*, which stores this same kind of information, albeit not formatted as three token sentences. When a database is intended to convey commonsense knowledge about everyday life or expert knowledge about an application area to an artificial intelligence system, we call the database a *knowledge base*. Knowledge bases range from general ones like Freebase, OpenCyc, WordNet, or Wikibase,<sup>1</sup> etc. to more specialized knowledge bases, like GeneOntology.<sup>2</sup> Representations for entities and relations can be learned by considering each triplet in a knowledge base as a training example and maximizing a training objective that captures their joint distribution (Bordes *et al.*, 2013a).

In addition to training data, we also need to define a model family to train. A common approach is to extend neural language models to model entities and relations. Neural language models learn a vector that provides a distributed representation of each word. They also learn about interactions between words, such as which word is likely to come after a sequence of words, by learning functions of these vectors. We can extend this approach to entities and relations by learning an embedding vector for each relation. In fact, the parallel between modeling

<sup>1</sup>Respectively available from these web sites: [freebase.com](http://freebase.com), [cyc.com/opencyc](http://cyc.com/opencyc), [wordnet.princeton.edu](http://wordnet.princeton.edu), [wikiba.se](http://wikiba.se)

<sup>2</sup>[geneontology.org](http://geneontology.org)

language and modeling knowledge encoded as relations is so close that researchers have trained representations of such entities by using *both* knowledge bases *and* natural language sentences (Bordes *et al.*, 2011, 2012; Wang *et al.*, 2014a) or combining data from multiple relational databases (Bordes *et al.*, 2013b). Many possibilities exist for the particular parametrization associated with such a model. Early work on learning about relations between entities (Paccanaro and Hinton, 2000) posited highly constrained parametric forms (“linear relational embeddings”), often using a different form of representation for the relation than for the entities. For example, Paccanaro and Hinton (2000) and Bordes *et al.* (2011) used vectors for entities and matrices for relations, with the idea that a relation acts like an operator on entities. Alternatively, relations can be considered as any other entity (Bordes *et al.*, 2012), allowing us to make statements about relations, but more flexibility is put in the machinery that combines them in order to model their joint distribution.

A practical short-term application of such models is *link prediction*: predicting missing arcs in the knowledge graph. This is a form of generalization to new facts, based on old facts. Most of the knowledge bases that currently exist have been constructed through manual labor, which tends to leave many and probably the majority of true relations absent from the knowledge base. See Wang *et al.* (2014b), Lin *et al.* (2015) and Garcia-Duran *et al.* (2015) for examples of such an application.

Evaluating the performance of a model on a link prediction task is difficult because we have only a dataset of positive examples (facts that are known to be true). If the model proposes a fact that is not in the dataset, we are unsure whether the model has made a mistake or discovered a new, previously unknown fact. The metrics are thus somewhat imprecise and are based on testing how the model ranks a held-out set of known true positive facts compared to other facts that are less likely to be true. A common way to construct interesting examples that are probably negative (facts that are probably false) is to begin with a true fact and create corrupted versions of that fact, for example by replacing one entity in the relation with a different entity selected at random. The popular precision at 10% metric counts how many times the model ranks a “correct” fact among the top 10% of all corrupted versions of that fact.

Another application of knowledge bases and distributed representations for them is *word-sense disambiguation* (Navigli and Velardi, 2005; Bordes *et al.*, 2012), which is the task of deciding which of the senses of a word is the appropriate one, in some context.

Eventually, knowledge of relations combined with a reasoning process and understanding of natural language could allow us to build a general question

answering system. A general question answering system must be able to process input information and remember important facts, organized in a way that enables it to retrieve and reason about them later. This remains a difficult open problem which can only be solved in restricted “toy” environments. Currently, the best approach to remembering and retrieving specific declarative facts is to use an explicit memory mechanism, as described in Sec. 10.12. Memory networks were first proposed to solve a toy question answering task (Weston *et al.*, 2014). Kumar *et al.* (2015) have proposed an extension that uses GRU recurrent nets to read the input into the memory and to produce the answer given the contents of the memory.

Deep learning has been applied to many other applications besides the ones described here, and will surely be applied to even more after this writing. It would be impossible to describe anything remotely resembling a comprehensive coverage of such a topic. This survey provides a representative sample of what is possible as of this writing.

This concludes Part II, which has described modern practices involving deep networks, comprising all of the most successful methods. Generally speaking, these methods involve using the gradient of a cost function to find the parameters of a model that approximates some desired function. With enough training data, this approach is extremely powerful. We now turn to Part III, in which we step into the territory of research—methods that are designed to work with less training data or to perform a greater variety of tasks, where the challenges are more difficult and not as close to being solved as the situations we have described so far.

## Part III

# Deep Learning Research

---

This part of the book describes the more ambitious and advanced approaches to deep learning, currently pursued by the research community.

In the previous parts of the book, we have shown how to solve supervised learning problems—how to learn to map one vector to another, given enough examples of the mapping.

Not all problems we might want to solve fall into this category. We may wish to generate new examples, or determine how likely some point is, or handle missing values and take advantage of a large set of unlabeled examples or examples from related tasks. A shortcoming of the current state of the art for industrial applications is that our learning algorithms require large amounts of supervised data to achieve good accuracy. In this part of the book, we discuss some of the speculative approaches to reducing the amount of labeled data necessary for existing models to work well and be applicable across a broader range of tasks. Accomplishing these goals usually requires some form of unsupervised or semi-supervised learning.

Many deep learning algorithms have been designed to tackle unsupervised learning problems, but none have truly solved the problem in the same way that deep learning has largely solved the supervised learning problem for a wide variety of tasks. In this part of the book, we describe the existing approaches to unsupervised learning and some of the popular thought about how we can make progress in this field.

A central cause of the difficulties with unsupervised learning is the high dimensionality of the random variables being modeled. This brings two distinct challenges: a statistical challenge and a computational challenge. The *statistical challenge* regards generalization: the number of configurations we may want to distinguish can grow exponentially with the number of dimensions of interest, and this quickly becomes much larger than the number of examples one can possibly have (or use with bounded computational resources). The *computational challenge* associated with high-dimensional distributions arises because many algorithms for learning or using a trained model (especially those based on estimating an explicit probability function) involve intractable computations that grow exponentially with the number of dimensions.

With probabilistic models, this computational challenge arises from the need to perform intractable inference or simply from the need to normalize the distribution.

- **Intractable inference:** inference is discussed mostly in Chapter 19. It regards the question of guessing the probable values of some variables  $a$ , given other variables  $b$ , with respect to a model that captures the joint

---

distribution between  $a$ ,  $b$  and  $c$ . In order to even compute such conditional probabilities one needs to sum over the values of the variables  $c$ , as well as compute a normalization constant which sums over the values of  $a$  and  $c$ .

- **Intractable normalization constants (the partition function):** the partition function is discussed mostly in Chapter 18. Normalizing constants of probability functions come up in inference (above) as well as in learning. Many probabilistic models involve such a normalizing constant. Unfortunately, learning such a model often requires computing the gradient of the logarithm of the partition function with respect to the model parameters. That computation is generally as intractable as computing the partition function itself. Monte Carlo Markov chain (MCMC) methods (Chapter 17) are often used to deal with the partition function (computing it or its gradient). Unfortunately, MCMC methods suffer when the modes of the model distribution are numerous and well-separated, especially in high-dimensional spaces (Sec. 17.5).

One way to confront these intractable computations is to approximate them, and many approaches have been proposed as discussed in this third part of the book. Another interesting way, also discussed here, would be to avoid these intractable computations altogether by design, and methods that do not require such computations are thus very appealing. Several generative models have been proposed in recent years, with that motivation. A wide variety of contemporary approaches to generative modeling are discussed in Chapter 20.

Part III is the most important for a researcher—someone who wants to understand the breadth of perspectives that have been brought to the field of deep learning, and push the field forward towards true artificial intelligence.

# Chapter 13

## Linear Factor Models

Many of the research frontiers in deep learning involve building a probabilistic model of the input,  $p_{\text{model}}(\mathbf{x})$ . Such a model can, in principle, use probabilistic inference to predict any of the variables in its environment given any of the other variables. Many of these models also have latent variables  $\mathbf{h}$ , with  $p_{\text{model}}(\mathbf{x}) = \mathbb{E}_{\mathbf{h}} p_{\text{model}}(\mathbf{x} | \mathbf{h})$ . These latent variables provide another means of representing the data. Distributed representations based on latent variables can obtain all of the advantages of representation learning that we have seen with deep feedforward and recurrent networks.

In this chapter, we describe some of the simplest probabilistic models with latent variables: linear factor models. These models are sometimes used as building blocks of mixture models (Hinton *et al.*, 1995a; Ghahramani and Hinton, 1996; Roweis *et al.*, 2002) or larger, deep probabilistic models (Tang *et al.*, 2012). They also show many of the basic approaches necessary to build generative models that the more advanced deep models will extend further.

A linear factor model is defined by the use of a stochastic, linear decoder function that generates  $\mathbf{x}$  by adding noise to a linear transformation of  $\mathbf{h}$ .

These models are interesting because they allow us to discover explanatory factors that have a simple joint distribution. The simplicity of using a linear decoder made these models some of the first latent variable models to be extensively studied.

A linear factor model describes the data generation process as follows. First, we sample the explanatory factors  $\mathbf{h}$  from a distribution

$$\mathbf{h} \sim p(\mathbf{h}), \quad (13.1)$$

where  $p(\mathbf{h})$  is a factorial distribution, with  $p(\mathbf{h}) = \prod_i p(h_i)$ , so that it is easy to

sample from. Next we sample the real-valued observable variables given the factors:

$$\mathbf{x} = \mathbf{W}\mathbf{h} + \mathbf{b} + \text{noise} \quad (13.2)$$

where the noise is typically Gaussian and diagonal (independent across dimensions). This is illustrated in Fig. 13.1.

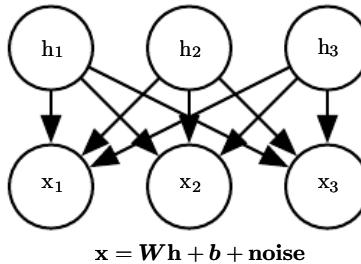


Figure 13.1: The directed graphical model describing the linear factor model family, in which we assume that an observed data vector  $\mathbf{x}$  is obtained by a linear combination of independent latent factors  $\mathbf{h}$ , plus some noise. Different models, such as probabilistic PCA, factor analysis or ICA, make different choices about the form of the noise and of the prior  $p(\mathbf{h})$ .

### 13.1 Probabilistic PCA and Factor Analysis

Probabilistic PCA (principal components analysis), factor analysis and other linear factor models are special cases of the above equations (13.1 and 13.2) and only differ in the choices made for the model's prior over latent variables  $\mathbf{h}$  before observing  $\mathbf{x}$  and noise distributions.

In *factor analysis* (Bartholomew, 1987; Basilevsky, 1994), the latent variable prior is just the unit variance Gaussian

$$\mathbf{h} \sim \mathcal{N}(\mathbf{h}; \mathbf{0}, \mathbf{I}) \quad (13.3)$$

while the observed variables  $x_i$  are assumed to be *conditionally independent*, given  $\mathbf{h}$ . Specifically, the noise is assumed to be drawn from a diagonal covariance Gaussian distribution, with covariance matrix  $\boldsymbol{\psi} = \text{diag}(\boldsymbol{\sigma}^2)$ , with  $\boldsymbol{\sigma}^2 = [\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2]^\top$  a vector of per-variable variances.

The role of the latent variables is thus to *capture the dependencies* between the different observed variables  $x_i$ . Indeed, it can easily be shown that  $\mathbf{x}$  is just a multivariate normal random variable, with

$$\mathbf{x} \sim \mathcal{N}(\mathbf{x}; \mathbf{b}, \mathbf{W}\mathbf{W}^\top + \boldsymbol{\psi}). \quad (13.4)$$

In order to cast PCA in a probabilistic framework, we can make a slight modification to the factor analysis model, making the conditional variances  $\sigma_i^2$  equal to each other. In that case the covariance of  $\mathbf{x}$  is just  $\mathbf{WW}^\top + \sigma^2\mathbf{I}$ , where  $\sigma^2$  is now a scalar. This yields the conditional distribution

$$\mathbf{x} \sim \mathcal{N}(\mathbf{x}; \mathbf{b}, \mathbf{WW}^\top + \sigma^2\mathbf{I}) \quad (13.5)$$

or equivalently

$$\mathbf{x} = \mathbf{W}\mathbf{h} + \mathbf{b} + \sigma\mathbf{z} \quad (13.6)$$

where  $\mathbf{z} \sim \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$  is Gaussian noise. Tipping and Bishop (1999) then show an iterative EM algorithm for estimating the parameters  $\mathbf{W}$  and  $\sigma^2$ .

This *probabilistic PCA* model takes advantage of the observation that most variations in the data can be captured by the latent variables  $\mathbf{h}$ , up to some small residual *reconstruction error*  $\sigma^2$ . As shown by Tipping and Bishop (1999), probabilistic PCA becomes PCA as  $\sigma \rightarrow 0$ . In that case, the conditional expected value of  $\mathbf{h}$  given  $\mathbf{x}$  becomes an orthogonal projection of  $\mathbf{x} - \mathbf{b}$  onto the space spanned by the  $d$  columns of  $\mathbf{W}$ , like in PCA.

As  $\sigma \rightarrow 0$ , the density model defined by probabilistic PCA becomes very sharp around these  $d$  dimensions spanned by the columns of  $\mathbf{W}$ . This can make the model assign very low likelihood to the data if the data does not actually cluster near a hyperplane.

## 13.2 Independent Component Analysis (ICA)

Independent component analysis (ICA) is among the oldest representation learning algorithms (Herault and Ans, 1984; Jutten and Herault, 1991; Comon, 1994; Hyvärinen, 1999; Hyvärinen *et al.*, 2001a; Hinton *et al.*, 2001; Teh *et al.*, 2003). It is an approach to modeling linear factors that seeks to separate an observed signal into many underlying signals that are scaled and added together to form the observed data. These signals are intended to be fully independent, rather than merely decorrelated from each other.<sup>1</sup>

Many different specific methodologies are referred to as ICA. The variant that is most similar to the other generative models we have described here is a variant (Pham *et al.*, 1992) that trains a fully parametric generative model. The prior distribution over the underlying factors,  $p(\mathbf{h})$ , must be fixed ahead of time by the user. The model then deterministically generates  $\mathbf{x} = \mathbf{W}\mathbf{h}$ . We can perform

---

<sup>1</sup>See Sec. 3.8 for a discussion of the difference between uncorrelated variables and independent variables.

a nonlinear change of variables (using Eq. 3.47) to determine  $p(\mathbf{x})$ . Learning the model then proceeds as usual, using maximum likelihood.

The motivation for this approach is that by choosing  $p(\mathbf{h})$  to be independent, we can recover underlying factors that are as close as possible to independent. This is commonly used, not to capture high-level abstract causal factors, but to recover low-level signals that have been mixed together. In this setting, each training example is one moment in time, each  $x_i$  is one sensor's observation of the mixed signals, and each  $h_i$  is one estimate of one of the original signals. For example, we might have  $n$  people speaking simultaneously. If we have  $n$  different microphones placed in different locations, ICA can detect the changes in the volume between each speaker as heard by each microphone, and separate the signals so that each  $h_i$  contains only one person speaking clearly. This is commonly used in neuroscience for electroencephalography, a technology for recording electrical signals originating in the brain. Many electrode sensors placed on the subject's head are used to measure many electrical signals coming from the body. The experimenter is typically only interested in signals from the brain, but signals from the subject's heart and eyes are strong enough to confound measurements taken at the subject's scalp. The signals arrive at the electrodes mixed together, so ICA is necessary to separate the electrical signature of the heart from the signals originating in the brain, and to separate signals in different brain regions from each other.

As mentioned before, many variants of ICA are possible. Some add some noise in the generation of  $\mathbf{x}$  rather than using a deterministic decoder. Most do not use the maximum likelihood criterion, but instead aim to make the elements of  $\mathbf{h} = \mathbf{W}^{-1}\mathbf{x}$  independent from each other. Many criteria that accomplish this goal are possible. Eq. 3.47 requires taking the determinant of  $\mathbf{W}$ , which can be an expensive and numerically unstable operation. Some variants of ICA avoid this problematic operation by constraining  $\mathbf{W}$  to be orthonormal.

All variants of ICA require that  $p(\mathbf{h})$  be non-Gaussian. This is because if  $p(\mathbf{h})$  is an independent prior with Gaussian components, then  $\mathbf{W}$  is not identifiable. We can obtain the same distribution over  $p(\mathbf{x})$  for many values of  $\mathbf{W}$ . This is very different from other linear factor models like probabilistic PCA and factor analysis, that often require  $p(\mathbf{h})$  to be Gaussian in order to make many operations on the model have closed form solutions. In the maximum likelihood approach where the user explicitly specifies the distribution, a typical choice is to use  $p(h_i) = \frac{d}{dh_i}\sigma(h_i)$ . Typical choices of these non-Gaussian distributions have larger peaks near 0 than does the Gaussian distribution, so we can also see most implementations of ICA as learning sparse features.

Many variants of ICA are not generative models in the sense that we use the phrase. In this book, a generative model either represents  $p(\mathbf{x})$  or can draw samples from it. Many variants of ICA only know how to transform between  $\mathbf{x}$  and  $\mathbf{h}$ , but do not have any way of representing  $p(\mathbf{h})$ , and thus do not impose a distribution over  $p(\mathbf{x})$ . For example, many ICA variants aim to increase the sample kurtosis of  $\mathbf{h} = \mathbf{W}^{-1}\mathbf{x}$ , because high kurtosis indicates that  $p(\mathbf{h})$  is non-Gaussian, but this is accomplished without explicitly representing  $p(\mathbf{h})$ . This is because ICA is more often used as an analysis tool for separating signals, rather than for generating data or estimating its density.

Just as PCA can be generalized to the nonlinear autoencoders described in Chapter 14, ICA can be generalized to a nonlinear generative model, in which we use a nonlinear function  $f$  to generate the observed data. See Hyvärinen and Pajunen (1999) for the initial work on nonlinear ICA and its successful use with ensemble learning by Roberts and Everson (2001) and Lappalainen *et al.* (2000). Another nonlinear extension of ICA is the approach of *nonlinear independent components estimation*, or NICE (Dinh *et al.*, 2014), which stacks a series of invertible transformations (encoder stages) that have the property that the determinant of the Jacobian of each transformation can be computed efficiently. This makes it possible to compute the likelihood exactly and, like ICA, attempts to transform the data into a space where it has a factorized marginal distribution, but is more likely to succeed thanks to the nonlinear encoder. Because the encoder is associated with a decoder that is its perfect inverse, it is straightforward to generate samples from the model (by first sampling from  $p(\mathbf{h})$  and then applying the decoder).

Another generalization of ICA is to learn groups of features, with statistical dependence allowed within a group but discouraged between groups (Hyvärinen and Hoyer, 1999; Hyvärinen *et al.*, 2001b). When the groups of related units are chosen to be non-overlapping, this is called *independent subspace analysis*. It is also possible to assign spatial coordinates to each hidden unit and form overlapping groups of spatially neighboring units. This encourages nearby units to learn similar features. When applied to natural images, this *topographic ICA* approach learns Gabor filters, such that neighboring features have similar orientation, location or frequency. Many different phase offsets of similar Gabor functions occur within each region, so that pooling over small regions yields translation invariance.

### 13.3 Slow Feature Analysis

*Slow feature analysis* (SFA) is a linear factor model that uses information from

time signals to learn invariant features (Wiskott and Sejnowski, 2002).

Slow feature analysis is motivated by a general principle called the slowness principle. The idea is that the important characteristics of scenes change very slowly compared to the individual measurements that make up a description of a scene. For example, in computer vision, individual pixel values can change very rapidly. If a zebra moves from left to right across the image, an individual pixel will rapidly change from black to white and back again as the zebra's stripes pass over the pixel. By comparison, the feature indicating whether a zebra is in the image will not change at all, and the feature describing the zebra's position will change slowly. We therefore may wish to regularize our model to learn features that change slowly over time.

The slowness principle predates slow feature analysis and has been applied to a wide variety of models (Hinton, 1989; Földiák, 1989; Mabahi *et al.*, 2009; Bergstra and Bengio, 2009). In general, we can apply the slowness principle to any differentiable model trained with gradient descent. The slowness principle may be introduced by adding a term to the cost function of the form

$$\lambda \sum_t L(f(\mathbf{x}^{(t+1)}), f(\mathbf{x}^{(t)})) \quad (13.7)$$

where  $\lambda$  is a hyperparameter determining the strength of the slowness regularization term,  $t$  is the index into a time sequence of examples,  $f$  is the feature extractor to be regularized, and  $L$  is a loss function measuring the distance between  $f(\mathbf{x}^{(t)})$  and  $f(\mathbf{x}^{(t+1)})$ . A common choice for  $L$  is the mean squared difference.

Slow feature analysis is a particularly efficient application of the slowness principle. It is efficient because it is applied to a linear feature extractor, and can thus be trained in closed form. Like some variants of ICA, SFA is not quite a generative model per se, in the sense that it defines a linear map between input space and feature space but does not define a prior over feature space and thus does not impose a distribution  $p(\mathbf{x})$  on input space.

The SFA algorithm (Wiskott and Sejnowski, 2002) consists of defining  $f(\mathbf{x}; \boldsymbol{\theta})$  to be a linear transformation, and solving the optimization problem

$$\min_{\boldsymbol{\theta}} \mathbb{E}_t (f(\mathbf{x}^{(t+1)})_i - f(\mathbf{x}^{(t)})_i)^2 \quad (13.8)$$

subject to the constraints

$$\mathbb{E}_t f(\mathbf{x}^{(t)})_i = 0 \quad (13.9)$$

and

$$\mathbb{E}_t [f(\mathbf{x}^{(t)})_i^2] = 1. \quad (13.10)$$

The constraint that the learned feature have zero mean is necessary to make the problem have a unique solution; otherwise we could add a constant to all feature values and obtain a different solution with equal value of the slowness objective. The constraint that the features have unit variance is necessary to prevent the pathological solution where all features collapse to 0. Like PCA, the SFA features are ordered, with the first feature being the slowest. To learn multiple features, we must also add the constraint

$$\forall i < j, \mathbb{E}_t [f(\mathbf{x}^{(t)})_i f(\mathbf{x}^{(t)})_j] = 0. \quad (13.11)$$

This specifies that the learned features must be linearly decorrelated from each other. Without this constraint, all of the learned features would simply capture the one slowest signal. One could imagine using other mechanisms, such as minimizing reconstruction error, to force the features to diversify, but this decorrelation mechanism admits a simple solution due to the linearity of SFA features. The SFA problem may be solved in closed form by a linear algebra package.

SFA is typically used to learn nonlinear features by applying a nonlinear basis expansion to  $\mathbf{x}$  before running SFA. For example, it is common to replace  $\mathbf{x}$  by the quadratic basis expansion, a vector containing elements  $x_i x_j$  for all  $i$  and  $j$ . Linear SFA modules may then be composed to learn deep nonlinear slow feature extractors by repeatedly learning a linear SFA feature extractor, applying a nonlinear basis expansion to its output, and then learning another linear SFA feature extractor on top of that expansion.

When trained on small spatial patches of videos of natural scenes, SFA with quadratic basis expansions learns features that share many characteristics with those of complex cells in V1 cortex (Berkes and Wiskott, 2005). When trained on videos of random motion within 3-D computer rendered environments, deep SFA learns features that share many characteristics with the features represented by neurons in rat brains that are used for navigation (Franzius *et al.*, 2007). SFA thus seems to be a reasonably biologically plausible model.

A major advantage of SFA is that it is possible to theoretically predict which features SFA will learn, even in the deep, nonlinear setting. To make such theoretical predictions, one must know about the dynamics of the environment in terms of configuration space (e.g., in the case of random motion in the 3-D rendered environment, the theoretical analysis proceeds from knowledge of the probability distribution over position and velocity of the camera). Given the knowledge of how the underlying factors actually change, it is possible to analytically solve for the optimal functions expressing these factors. In practice, experiments with deep SFA applied to simulated data seem to recover the theoretically predicted functions.

This is in comparison to other learning algorithms where the cost function depends highly on specific pixel values, making it much more difficult to determine what features the model will learn.

Deep SFA has also been used to learn features for object recognition and pose estimation (Franzius *et al.*, 2008). So far, the slowness principle has not become the basis for any state of the art applications. It is unclear what factor has limited its performance. We speculate that perhaps the slowness prior is too strong, and that, rather than imposing a prior that features should be approximately constant, it would be better to impose a prior that features should be easy to predict from one time step to the next. The position of an object is a useful feature regardless of whether the object’s velocity is high or low, but the slowness principle encourages the model to ignore the position of objects that have high velocity.

## 13.4 Sparse Coding

*Sparse coding* (Olshausen and Field, 1996) is a linear factor model that has been heavily studied as an unsupervised feature learning and feature extraction mechanism. Strictly speaking, the term “sparse coding” refers to the process of inferring the value of  $\mathbf{h}$  in this model, while “sparse modeling” refers to the process of designing and learning the model, but the term “sparse coding” is often used to refer to both.

Like most other linear factor models, it uses a linear decoder plus noise to obtain reconstructions of  $\mathbf{x}$ , as specified in Eq. 13.2. More specifically, sparse coding models typically assume that the linear factors have Gaussian noise with isotropic precision  $\beta$ :

$$p(\mathbf{x} \mid \mathbf{h}) = \mathcal{N}(\mathbf{x}; \mathbf{W}\mathbf{h} + \mathbf{b}, \frac{1}{\beta}\mathbf{I}). \quad (13.12)$$

The distribution  $p(\mathbf{h})$  is chosen to be one with sharp peaks near 0 (Olshausen and Field, 1996). Common choices include factorized Laplace, Cauchy or factorized Student- $t$  distributions. For example, the Laplace prior parametrized in terms of the sparsity penalty coefficient  $\lambda$  is given by

$$p(h_i) = \text{Laplace}(h_i; 0, \frac{2}{\lambda}) = \frac{\lambda}{4} e^{-\frac{1}{2}\lambda|h_i|} \quad (13.13)$$

and the Student- $t$  prior by

$$p(h_i) \propto \frac{1}{(1 + \frac{h_i^2}{\nu})^{\frac{\nu+1}{2}}}. \quad (13.14)$$

Training sparse coding with maximum likelihood is intractable. Instead, the training alternates between encoding the data and training the decoder to better reconstruct the data given the encoding. This approach will be justified further as a principled approximation to maximum likelihood later, in Sec. 19.3.

For models such as PCA, we have seen the use of a parametric encoder function that predicts  $\mathbf{h}$  and consists only of multiplication by a weight matrix. The encoder that we use with sparse coding is not a parametric encoder. Instead, the encoder is an optimization algorithm, that solves an optimization problem in which we seek the single most likely code value:

$$\mathbf{h}^* = f(\mathbf{x}) = \arg \max_{\mathbf{h}} p(\mathbf{h} | \mathbf{x}). \quad (13.15)$$

When combined with Eq. 13.13 and Eq. 13.12, this yields the following optimization problem:

$$\arg \max_{\mathbf{h}} p(\mathbf{h} | \mathbf{x}) \quad (13.16)$$

$$= \arg \max_{\mathbf{h}} \log p(\mathbf{h} | \mathbf{x}) \quad (13.17)$$

$$= \arg \min_{\mathbf{h}} \lambda \|\mathbf{h}\|_1 + \beta \|\mathbf{x} - \mathbf{W}\mathbf{h}\|_2^2, \quad (13.18)$$

where we have dropped terms not depending on  $\mathbf{h}$  and divided by positive scaling factors to simplify the equation.

Due to the imposition of an  $L^1$  norm on  $\mathbf{h}$ , that this procedure will yield a sparse  $\mathbf{h}^*$  (See Sec. 7.1.2).

To train the model rather than just perform inference, we alternate between minimization with respect to  $\mathbf{h}$  and minimization with respect to  $\mathbf{W}$ . In this presentation, we treat  $\beta$  as a hyperparameter. Typically it is set to 1 because its role in this optimization problem is shared with  $\lambda$  and there is no need for both hyperparameters. In principle, we could also treat  $\beta$  as a parameter of the model and learn it. Our presentation here has discarded some terms that do not depend on  $\mathbf{h}$  but do depend on  $\beta$ . To learn  $\beta$ , these terms must be included, or  $\beta$  will collapse to 0.

Not all approaches to sparse coding explicitly build a  $p(\mathbf{h})$  and a  $p(\mathbf{x} | \mathbf{h})$ . Often we are just interested in learning a dictionary of features with activation values that will often be zero when extracted using this inference procedure.

If we sample  $\mathbf{h}$  from a Laplace prior, it is in fact a zero probability event for an element of  $\mathbf{h}$  to actually be zero. The generative model itself is not especially sparse, only the feature extractor is. Goodfellow *et al.* (2013d) describe approximate

inference in a different model family, the spike and slab sparse coding model, for which samples from the prior usually contain true zeros.

The sparse coding approach combined with the use of the non-parametric encoder can in principle minimize the combination of reconstruction error and log-prior better than any specific parametric encoder. Another advantage is that there is no generalization error to the encoder. A parametric encoder must learn how to map  $\mathbf{x}$  to  $\mathbf{h}$  in a way that generalizes. For unusual  $\mathbf{x}$  that do not resemble the training data, a learned, parametric encoder may fail to find an  $\mathbf{h}$  that results in accurate reconstruction or a sparse code. For the vast majority of formulations of sparse coding models, where the inference problem is convex, the optimization procedure will always find the optimal code (unless degenerate cases such as replicated weight vectors occur). Obviously, the sparsity and reconstruction costs can still rise on unfamiliar points, but this is due to generalization error in the decoder weights, rather than generalization error in the encoder. The lack of generalization error in sparse coding's optimization-based encoding process may result in better generalization when sparse coding is used as a feature extractor for a classifier than when a parametric function is used to predict the code. [Coates and Ng \(2011\)](#) demonstrated that sparse coding features generalize better for object recognition tasks than the features of a related model based on a parametric encoder, the linear-sigmoid autoencoder. Inspired by their work, [Goodfellow et al. \(2013d\)](#) showed that a variant of sparse coding generalizes better than other feature extractors in the regime where extremely few labels are available (twenty or fewer labels per class).

The primary disadvantage of the non-parametric encoder is that it requires greater time to compute  $\mathbf{h}$  given  $\mathbf{x}$  because the non-parametric approach requires running an iterative algorithm. The parametric autoencoder approach, developed in Chapter 14, uses only a fixed number of layers, often only one. Another disadvantage is that it is not straight-forward to back-propagate through the non-parametric encoder, which makes it difficult to pretrain a sparse coding model with an unsupervised criterion and then fine-tune it using a supervised criterion. Modified versions of sparse coding that permit approximate derivatives do exist but are not widely used ([Bagnell and Bradley, 2009](#)).

Sparse coding, like other linear factor models, often produces poor samples, as shown in Fig. 13.2. This happens even when the model is able to reconstruct the data well and provide useful features for a classifier. The reason is that each individual feature may be learned well, but the factorial prior on the hidden code results in the model including random subsets of all of the features in each generated sample. This motivates the development of deeper models that can impose a non-

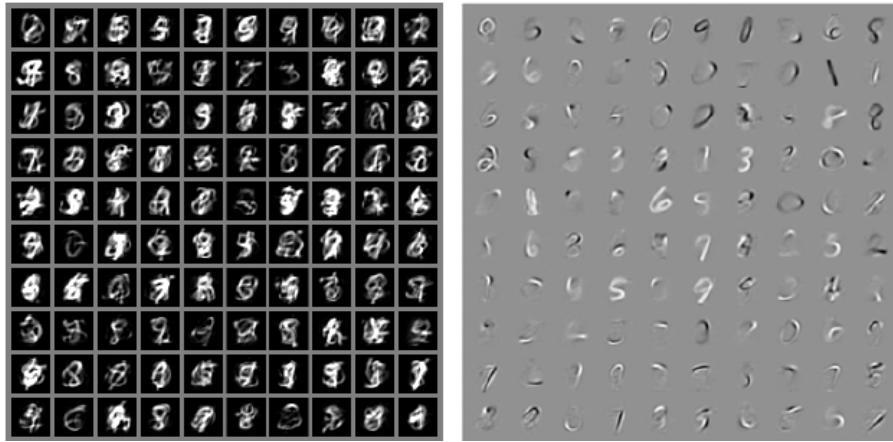


Figure 13.2: Example samples and weights from a spike and slab sparse coding model trained on the MNIST dataset. (*Left*) The samples from the model do not resemble the training examples. At first glance, one might assume the model is poorly fit. (*Right*) The weight vectors of the model have learned to represent penstrokes and sometimes complete digits. The model has thus learned useful features. The problem is that the factorial prior over features results in random subsets of features being combined. Few such subsets are appropriate to form a recognizable MNIST digit. This motivates the development of generative models that have more powerful distributions over their latent codes. Figure reproduced with permission from [Goodfellow et al. \(2013d\)](#).

factorial distribution on the deepest code layer, as well as the development of more sophisticated shallow models.

## 13.5 Manifold Interpretation of PCA

Linear factor models including PCA and factor analysis can be interpreted as learning a manifold ([Hinton et al., 1997](#)). We can view probabilistic PCA as defining a thin pancake-shaped region of high probability—a Gaussian distribution that is very narrow along some axes, just as a pancake is very flat along its vertical axis, but is elongated along other axes, just as a pancake is wide along its horizontal axes. This is illustrated in Fig. 13.3. PCA can be interpreted as aligning this pancake with a linear manifold in a higher-dimensional space. This interpretation applies not just to traditional PCA but also to any linear autoencoder that learns matrices  $\mathbf{W}$  and  $\mathbf{V}$  with the goal of making the reconstruction of  $\mathbf{x}$  lie as close to  $\mathbf{x}$  as possible,

Let the encoder be

$$\mathbf{h} = f(\mathbf{x}) = \mathbf{W}^\top(\mathbf{x} - \boldsymbol{\mu}). \quad (13.19)$$

The encoder computes a low-dimensional representation of  $h$ . With the autoencoder view, we have a decoder computing the reconstruction

$$\hat{x} = g(\mathbf{h}) = \mathbf{b} + \mathbf{V}\mathbf{h}. \quad (13.20)$$

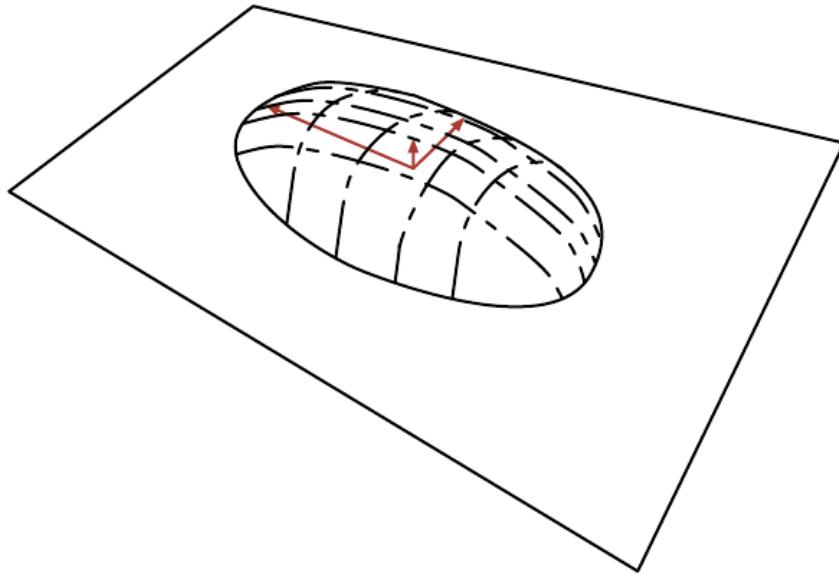


Figure 13.3: Flat Gaussian capturing probability concentration near a low-dimensional manifold. The figure shows the upper half of the “pancake” above the “manifold plane” which goes through its middle. The variance in the direction orthogonal to the manifold is very small (arrow pointing out of plane) and can be considered like “noise,” while the other variances are large (arrows in the plane) and correspond to “signal,” and a coordinate system for the reduced-dimension data.

The choices of linear encoder and decoder that minimize reconstruction error

$$\mathbb{E}[||\mathbf{x} - \hat{\mathbf{x}}||^2] \quad (13.21)$$

correspond to  $\mathbf{V} = \mathbf{W}$ ,  $\boldsymbol{\mu} = \mathbf{b} = \mathbb{E}[\mathbf{x}]$  and the columns of  $\mathbf{W}$  form an orthonormal basis which spans the same subspace as the principal eigenvectors of the covariance matrix

$$\mathbf{C} = \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top]. \quad (13.22)$$

In the case of PCA, the columns of  $\mathbf{W}$  are these eigenvectors, ordered by the magnitude of the corresponding eigenvalues (which are all real and non-negative).

One can also show that eigenvalue  $\lambda_i$  of  $\mathbf{C}$  corresponds to the variance of  $\mathbf{x}$  in the direction of eigenvector  $\mathbf{v}^{(i)}$ . If  $\mathbf{x} \in \mathbb{R}^D$  and  $\mathbf{h} \in \mathbb{R}^d$  with  $d < D$ , then the

optimal reconstruction error (choosing  $\mu$ ,  $b$ ,  $V$  and  $W$  as above) is

$$\min \mathbb{E}[||\mathbf{x} - \hat{\mathbf{x}}||^2] = \sum_{i=d+1}^D \lambda_i. \quad (13.23)$$

Hence, if the covariance has rank  $d$ , the eigenvalues  $\lambda_{d+1}$  to  $\lambda_D$  are 0 and reconstruction error is 0.

Furthermore, one can also show that the above solution can be obtained by maximizing the variances of the elements of  $\mathbf{h}$ , under orthonormal  $\mathbf{W}$ , instead of minimizing reconstruction error.

Linear factor models are some of the simplest generative models and some of the simplest models that learn a representation of data. Much as linear classifiers and linear regression models may be extended to deep feedforward networks, these linear factor models may be extended to autoencoder networks and deep probabilistic models that perform the same tasks but with a much more powerful and flexible model family.

# Chapter 14

## Autoencoders

An *autoencoder* is a neural network that is trained to attempt to copy its input to its output. Internally, it has a hidden layer  $\mathbf{h}$  that describes a *code* used to represent the input. The network may be viewed as consisting of two parts: an encoder function  $\mathbf{h} = f(\mathbf{x})$  and a decoder that produces a reconstruction  $\mathbf{r} = g(\mathbf{h})$ . This architecture is presented in Fig. 14.1. If an autoencoder succeeds in simply learning to set  $g(f(\mathbf{x})) = \mathbf{x}$  everywhere, then it is not especially useful. Instead, autoencoders are designed to be unable to learn to copy perfectly. Usually they are restricted in ways that allow them to copy only approximately, and to copy only input that resembles the training data. Because the model is forced to prioritize which aspects of the input should be copied, it often learns useful properties of the data.

Modern autoencoders have generalized the idea of an encoder and a decoder beyond deterministic functions to stochastic mappings  $p_{\text{encoder}}(\mathbf{h} | \mathbf{x})$  and  $p_{\text{decoder}}(\mathbf{x} | \mathbf{h})$ .

The idea of autoencoders has been part of the historical landscape of neural networks for decades (LeCun, 1987; Bourlard and Kamp, 1988; Hinton and Zemel, 1994). Traditionally, autoencoders were used for dimensionality reduction or feature learning. Recently, theoretical connections between autoencoders and latent variable models have brought autoencoders to the forefront of generative modeling, as we will see in Chapter 20. Autoencoders may be thought of as being a special case of feedforward networks, and may be trained with all of the same techniques, typically minibatch gradient descent following gradients computed by back-propagation. Unlike general feedforward networks, autoencoders may also be trained using *recirculation* (Hinton and McClelland, 1988), a learning algorithm based on comparing the activations of the network on the original input

to the activations on the reconstructed input. Recirculation is regarded as more biologically plausible than back-propagation, but is rarely used for machine learning applications.

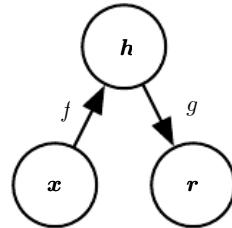


Figure 14.1: The general structure of an autoencoder, mapping an input  $\mathbf{x}$  to an output (called reconstruction)  $\mathbf{r}$  through an internal representation or code  $\mathbf{h}$ . The autoencoder has two components: the encoder  $f$  (mapping  $\mathbf{x}$  to  $\mathbf{h}$ ) and the decoder  $g$  (mapping  $\mathbf{h}$  to  $\mathbf{r}$ ).

## 14.1 Undercomplete Autoencoders

Copying the input to the output may sound useless, but we are typically not interested in the output of the decoder. Instead, we hope that training the autoencoder to perform the input copying task will result in  $\mathbf{h}$  taking on useful properties.

One way to obtain useful features from the autoencoder is to constrain  $\mathbf{h}$  to have smaller dimension than  $\mathbf{x}$ . An autoencoder whose code dimension is less than the input dimension is called *undercomplete*. Learning an undercomplete representation forces the autoencoder to capture the most salient features of the training data.

The learning process is described simply as minimizing a loss function

$$L(\mathbf{x}, g(f(\mathbf{x}))) \tag{14.1}$$

where  $L$  is a loss function penalizing  $g(f(\mathbf{x}))$  for being dissimilar from  $\mathbf{x}$ , such as the mean squared error.

When the decoder is linear and  $L$  is the mean squared error, an undercomplete autoencoder learns to span the same subspace as PCA. In this case, an autoencoder trained to perform the copying task has learned the principal subspace of the training data as a side-effect.

Autoencoders with nonlinear encoder functions  $f$  and nonlinear decoder functions  $g$  can thus learn a more powerful nonlinear generalization of PCA. Unfortu-

nately, if the encoder and decoder are allowed too much capacity, the autoencoder can learn to perform the copying task without extracting useful information about the distribution of the data. Theoretically, one could imagine that an autoencoder with a one-dimensional code but a very powerful nonlinear encoder could learn to represent each training example  $\mathbf{x}^{(i)}$  with the code  $i$ . The decoder could learn to map these integer indices back to the values of specific training examples. This specific scenario does not occur in practice, but it illustrates clearly that an autoencoder trained to perform the copying task can fail to learn anything useful about the dataset if the capacity of the autoencoder is allowed to become too great.

## 14.2 Regularized Autoencoders

Undercomplete autoencoders, with code dimension less than the input dimension, can learn the most salient features of the data distribution. We have seen that these autoencoders fail to learn anything useful if the encoder and decoder are given too much capacity.

A similar problem occurs if the hidden code is allowed to have dimension equal to the input, and in the *overcomplete* case in which the hidden code has dimension greater than the input. In these cases, even a linear encoder and linear decoder can learn to copy the input to the output without learning anything useful about the data distribution.

Ideally, one could train any architecture of autoencoder successfully, choosing the code dimension and the capacity of the encoder and decoder based on the complexity of distribution to be modeled. Regularized autoencoders provide the ability to do so. Rather than limiting the model capacity by keeping the encoder and decoder shallow and the code size small, regularized autoencoders use a loss function that encourages the model to have other properties besides the ability to copy its input to its output. These other properties include sparsity of the representation, smallness of the derivative of the representation, and robustness to noise or to missing inputs. A regularized autoencoder can be nonlinear and overcomplete but still learn something useful about the data distribution even if the model capacity is great enough to learn a trivial identity function.

In addition to the methods described here which are most naturally interpreted as regularized autoencoders, nearly any generative model with latent variables and equipped with an inference procedure (for computing latent representations given input) may be viewed as a particular form of autoencoder. Two generative modeling approaches that emphasize this connection with autoencoders are the descendants of the Helmholtz machine (Hinton *et al.*, 1995b), such as the variational

autoencoder (Sec. 20.10.3) and the generative stochastic networks (Sec. 20.12). These models naturally learn high-capacity, overcomplete encodings of the input and do not require regularization for these encodings to be useful. Their encodings are naturally useful because the models were trained to approximately maximize the probability of the training data rather than to copy the input to the output.

### 14.2.1 Sparse Autoencoders

A sparse autoencoder is simply an autoencoder whose training criterion involves a sparsity penalty  $\Omega(\mathbf{h})$  on the code layer  $\mathbf{h}$ , in addition to the reconstruction error:

$$L(\mathbf{x}, g(f(\mathbf{x}))) + \Omega(\mathbf{h}) \quad (14.2)$$

where  $g(\mathbf{h})$  is the decoder output and typically we have  $\mathbf{h} = f(\mathbf{x})$ , the encoder output.

Sparse autoencoders are typically used to learn features for another task such as classification. An autoencoder that has been regularized to be sparse must respond to unique statistical features of the dataset it has been trained on, rather than simply acting as an identity function. In this way, training to perform the copying task with a sparsity penalty can yield a model that has learned useful features as a byproduct.

We can think of the penalty  $\Omega(\mathbf{h})$  simply as a regularizer term added to a feedforward network whose primary task is to copy the input to the output (unsupervised learning objective) and possibly also perform some supervised task (with a supervised learning objective) that depends on these sparse features. Unlike other regularizers such as weight decay, there is not a straightforward Bayesian interpretation to this regularizer. As described in Sec. 5.6.1, training with weight decay and other regularization penalties can be interpreted as a MAP approximation to Bayesian inference, with the added regularizing penalty corresponding to a prior probability distribution over the model parameters. In this view, regularized maximum likelihood corresponds to maximizing  $p(\boldsymbol{\theta} | \mathbf{x})$ , which is equivalent to maximizing  $\log p(\mathbf{x} | \boldsymbol{\theta}) + \log p(\boldsymbol{\theta})$ . The  $\log p(\mathbf{x} | \boldsymbol{\theta})$  term is the usual data log-likelihood term and the  $\log p(\boldsymbol{\theta})$  term, the log-prior over parameters, incorporates the preference over particular values of  $\boldsymbol{\theta}$ . This view was described in Sec. 5.6. Regularized autoencoders defy such an interpretation because the regularizer depends on the data and is therefore by definition not a prior in the formal sense of the word. We can still think of these regularization terms as implicitly expressing a preference over functions.

Rather than thinking of the sparsity penalty as a regularizer for the copying task, we can think of the entire sparse autoencoder framework as approximating

maximum likelihood training of a generative model that has latent variables. Suppose we have a model with visible variables  $\mathbf{x}$  and latent variables  $\mathbf{h}$ , with an explicit joint distribution  $p_{\text{model}}(\mathbf{x}, \mathbf{h}) = p_{\text{model}}(\mathbf{h})p_{\text{model}}(\mathbf{x} | \mathbf{h})$ . We refer to  $p_{\text{model}}(\mathbf{h})$  as the model’s prior distribution over the latent variables, representing the model’s beliefs prior to seeing  $\mathbf{x}$ . This is different from the way we have previously used the word “prior,” to refer to the distribution  $p(\boldsymbol{\theta})$  encoding our beliefs about the model’s parameters before we have seen the training data. The log-likelihood can be decomposed as

$$\log p_{\text{model}}(\mathbf{x}) = \log \sum_{\mathbf{h}} p_{\text{model}}(\mathbf{h}, \mathbf{x}). \quad (14.3)$$

We can think of the autoencoder as approximating this sum with a point estimate for just one highly likely value for  $\mathbf{h}$ . This is similar to the sparse coding generative model (Sec. 13.4), but with  $\mathbf{h}$  being the output of the parametric encoder rather than the result of an optimization that infers the most likely  $\mathbf{h}$ . From this point of view, with this chosen  $\mathbf{h}$ , we are maximizing

$$\log p_{\text{model}}(\mathbf{h}, \mathbf{x}) = \log p_{\text{model}}(\mathbf{h}) + \log p_{\text{model}}(\mathbf{x} | \mathbf{h}). \quad (14.4)$$

The  $\log p_{\text{model}}(\mathbf{h})$  term can be sparsity-inducing. For example, the Laplace prior,

$$p_{\text{model}}(h_i) = \frac{\lambda}{2} e^{-\lambda|h_i|}, \quad (14.5)$$

corresponds to an absolute value sparsity penalty. Expressing the log-prior as an absolute value penalty, we obtain

$$\Omega(\mathbf{h}) = \lambda \sum_i |h_i| \quad (14.6)$$

$$-\log p_{\text{model}}(\mathbf{h}) = \sum_i \left( \lambda |h_i| - \log \frac{\lambda}{2} \right) = \Omega(\mathbf{h}) + \text{const} \quad (14.7)$$

where the constant term depends only on  $\lambda$  and not  $\mathbf{h}$ . We typically treat  $\lambda$  as a hyperparameter and discard the constant term since it does not affect the parameter learning. Other priors such as the Student- $t$  prior can also induce sparsity. From this point of view of sparsity as resulting from the effect of  $p_{\text{model}}(\mathbf{h})$  on approximate maximum likelihood learning, the sparsity penalty is not a regularization term at all. It is just a consequence of the model’s distribution over its latent variables. This view provides a different motivation for training an autoencoder: it is a way of approximately training a generative model. It also provides a different reason for

why the features learned by the autoencoder are useful: they describe the latent variables that explain the input.

Early work on sparse autoencoders ([Ranzato et al., 2007a, 2008](#)) explored various forms of sparsity and proposed a connection between the sparsity penalty and the  $\log Z$  term that arises when applying maximum likelihood to an undirected probabilistic model  $p(\mathbf{x}) = \frac{1}{Z} \tilde{p}(\mathbf{x})$ . The idea is that minimizing  $\log Z$  prevents a probabilistic model from having high probability everywhere, and imposing sparsity on an autoencoder prevents the autoencoder from having low reconstruction error everywhere. In this case, the connection is on the level of an intuitive understanding of a general mechanism rather than a mathematical correspondence. The interpretation of the sparsity penalty as corresponding to  $\log p_{\text{model}}(\mathbf{h})$  in a directed model  $p_{\text{model}}(\mathbf{h})p_{\text{model}}(\mathbf{x} | \mathbf{h})$  is more mathematically straightforward.

One way to achieve *actual zeros* in  $\mathbf{h}$  for sparse (and denoising) autoencoders was introduced in [Glorot et al. \(2011b\)](#). The idea is to use rectified linear units to produce the code layer. With a prior that actually pushes the representations to zero (like the absolute value penalty), one can thus indirectly control the average number of zeros in the representation.

### 14.2.2 Denoising Autoencoders

Rather than adding a penalty  $\Omega$  to the cost function, we can obtain an autoencoder that learns something useful by changing the reconstruction error term of the cost function.

Traditionally, autoencoders minimize some function

$$L(\mathbf{x}, g(f(\mathbf{x}))) \tag{14.8}$$

where  $L$  is a loss function penalizing  $g(f(\mathbf{x}))$  for being dissimilar from  $\mathbf{x}$ , such as the  $L^2$  norm of their difference. This encourages  $g \circ f$  to learn to be merely an identity function if they have the capacity to do so.

A *denoising autoencoder* or *DAE* instead minimizes

$$L(\mathbf{x}, g(f(\tilde{\mathbf{x}}))), \tag{14.9}$$

where  $\tilde{\mathbf{x}}$  is a copy of  $\mathbf{x}$  that has been corrupted by some form of noise. Denoising autoencoders must therefore undo this corruption rather than simply copying their input.

Denoising training forces  $f$  and  $g$  to implicitly learn the structure of  $p_{\text{data}}(\mathbf{x})$ , as shown by [Alain and Bengio \(2013\)](#) and [Bengio et al. \(2013c\)](#). Denoising

autoencoders thus provide yet another example of how useful properties can emerge as a byproduct of minimizing reconstruction error. They are also an example of how overcomplete, high-capacity models may be used as autoencoders so long as care is taken to prevent them from learning the identity function. Denoising autoencoders are presented in more detail in Sec. 14.5.

### 14.2.3 Regularizing by Penalizing Derivatives

Another strategy for regularizing an autoencoder is to use a penalty  $\Omega$  as in sparse autoencoders,

$$L(\mathbf{x}, g(f(\mathbf{x}))) + \Omega(\mathbf{h}, \mathbf{x}), \quad (14.10)$$

but with a different form of  $\Omega$ :

$$\Omega(\mathbf{h}, \mathbf{x}) = \lambda \sum_i \|\nabla_{\mathbf{x}} h_i\|^2. \quad (14.11)$$

This forces the model to learn a function that does not change much when  $\mathbf{x}$  changes slightly. Because this penalty is applied only at training examples, it forces the autoencoder to learn features that capture information about the training distribution.

An autoencoder regularized in this way is called a *contractive autoencoder* or *CAE*. This approach has theoretical connections to denoising autoencoders, manifold learning and probabilistic modeling. The CAE is described in more detail in Sec. 14.7.

## 14.3 Representational Power, Layer Size and Depth

Autoencoders are often trained with only a single layer encoder and a single layer decoder. However, this is not a requirement. In fact, using deep encoders and decoders offers many advantages.

Recall from Sec. 6.4.1 that there are many advantages to depth in a feedforward network. Because autoencoders are feedforward networks, these advantages also apply to autoencoders. Moreover, the encoder is itself a feedforward network as is the decoder, so each of these components of the autoencoder can individually benefit from depth.

One major advantage of non-trivial depth is that the universal approximator theorem guarantees that a feedforward neural network with at least one hidden layer can represent an approximation of any function (within a broad class) to an

arbitrary degree of accuracy, provided that it has enough hidden units. This means that an autoencoder with a single hidden layer is able to represent the identity function along the domain of the data arbitrarily well. However, the mapping from input to code is shallow. This means that we are not able to enforce arbitrary constraints, such as that the code should be sparse. A deep autoencoder, with at least one additional hidden layer inside the encoder itself, can approximate any mapping from input to code arbitrarily well, given enough hidden units.

Depth can exponentially reduce the computational cost of representing some functions. Depth can also exponentially decrease the amount of training data needed to learn some functions. See Sec. 6.4.1 for a review of the advantages of depth in feedforward networks.

Experimentally, deep autoencoders yield much better compression than corresponding shallow or linear autoencoders (Hinton and Salakhutdinov, 2006).

A common strategy for training a deep autoencoder is to greedily pretrain the deep architecture by training a stack of shallow autoencoders, so we often encounter shallow autoencoders, even when the ultimate goal is to train a deep autoencoder.

## 14.4 Stochastic Encoders and Decoders

Autoencoders are just feedforward networks. The same loss functions and output unit types that can be used for traditional feedforward networks are also used for autoencoders.

As described in Sec. 6.2.2.4, a general strategy for designing the output units and the loss function of a feedforward network is to define an output distribution  $p(\mathbf{y} | \mathbf{x})$  and minimize the negative log-likelihood  $-\log p(\mathbf{y} | \mathbf{x})$ . In that setting,  $\mathbf{y}$  was a vector of targets, such as class labels.

In the case of an autoencoder,  $\mathbf{x}$  is now the target as well as the input. However, we can still apply the same machinery as before. Given a hidden code  $\mathbf{h}$ , we may think of the decoder as providing a conditional distribution  $p_{\text{decoder}}(\mathbf{x} | \mathbf{h})$ . We may then train the autoencoder by minimizing  $-\log p_{\text{decoder}}(\mathbf{x} | \mathbf{h})$ . The exact form of this loss function will change depending on the form of  $p_{\text{decoder}}$ . As with traditional feedforward networks, we usually use linear output units to parametrize the mean of a Gaussian distribution if  $\mathbf{x}$  is real-valued. In that case, the negative log-likelihood yields a mean squared error criterion. Similarly, binary  $\mathbf{x}$  values correspond to a Bernoulli distribution whose parameters are given by a sigmoid output unit, discrete  $\mathbf{x}$  values correspond to a softmax distribution, and so on.

Typically, the output variables are treated as being conditionally independent given  $\mathbf{h}$  so that this probability distribution is inexpensive to evaluate, but some techniques such as mixture density outputs allow tractable modeling of outputs with correlations.

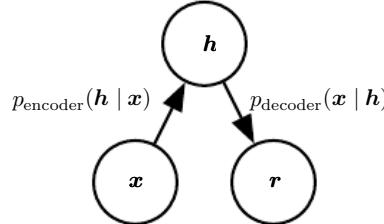


Figure 14.2: The structure of a stochastic autoencoder, in which both the encoder and the decoder are not simple functions but instead involve some noise injection, meaning that their output can be seen as sampled from a distribution,  $p_{\text{encoder}}(\mathbf{h} | \mathbf{x})$  for the encoder and  $p_{\text{decoder}}(\mathbf{x} | \mathbf{h})$  for the decoder.

To make a more radical departure from the feedforward networks we have seen previously, we can also generalize the notion of an *encoding function*  $f(\mathbf{x})$  to an *encoding distribution*  $p_{\text{encoder}}(\mathbf{h} | \mathbf{x})$ , as illustrated in Fig. 14.2.

Any latent variable model  $p_{\text{model}}(\mathbf{h}, \mathbf{x})$  defines a stochastic encoder

$$p_{\text{encoder}}(\mathbf{h} | \mathbf{x}) = p_{\text{model}}(\mathbf{h} | \mathbf{x}) \quad (14.12)$$

and a stochastic decoder

$$p_{\text{decoder}}(\mathbf{x} | \mathbf{h}) = p_{\text{model}}(\mathbf{x} | \mathbf{h}). \quad (14.13)$$

In general, the encoder and decoder distributions are not necessarily conditional distributions compatible with a unique joint distribution  $p_{\text{model}}(\mathbf{x}, \mathbf{h})$ . Alain *et al.* (2015) showed that training the encoder and decoder as a denoising autoencoder will tend to make them compatible asymptotically (with enough capacity and examples).

## 14.5 Denoising Autoencoders

The *denoising autoencoder* (DAE) is an autoencoder that receives a corrupted data point as input and is trained to predict the original, uncorrupted data point as its output.

The DAE training procedure is illustrated in Fig. 14.3. We introduce a corruption process  $C(\tilde{\mathbf{x}} | \mathbf{x})$  which represents a conditional distribution over

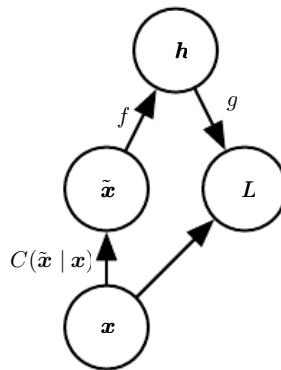


Figure 14.3: The computational graph of the cost function for a denoising autoencoder, which is trained to reconstruct the clean data point  $\mathbf{x}$  from its corrupted version  $\tilde{\mathbf{x}}$ . This is accomplished by minimizing the loss  $L = -\log p_{\text{decoder}}(\mathbf{x} \mid \mathbf{h} = f(\tilde{\mathbf{x}}))$ , where  $\tilde{\mathbf{x}}$  is a corrupted version of the data example  $\mathbf{x}$ , obtained through a given corruption process  $C(\tilde{\mathbf{x}} \mid \mathbf{x})$ . Typically the distribution  $p_{\text{decoder}}$  is a factorial distribution whose mean parameters are emitted by a feedforward network  $g$ .

corrupted samples  $\tilde{\mathbf{x}}$ , given a data sample  $\mathbf{x}$ . The autoencoder then learns a *reconstruction distribution*  $p_{\text{reconstruct}}(\mathbf{x} \mid \tilde{\mathbf{x}})$  estimated from training pairs  $(\mathbf{x}, \tilde{\mathbf{x}})$ , as follows:

1. Sample a training example  $\mathbf{x}$  from the training data.
2. Sample a corrupted version  $\tilde{\mathbf{x}}$  from  $C(\tilde{\mathbf{x}} \mid \mathbf{x} = \mathbf{x})$ .
3. Use  $(\mathbf{x}, \tilde{\mathbf{x}})$  as a training example for estimating the autoencoder reconstruction distribution  $p_{\text{reconstruct}}(\mathbf{x} \mid \tilde{\mathbf{x}}) = p_{\text{decoder}}(\mathbf{x} \mid \mathbf{h})$  with  $\mathbf{h}$  the output of encoder  $f(\tilde{\mathbf{x}})$  and  $p_{\text{decoder}}$  typically defined by a decoder  $g(\mathbf{h})$ .

Typically we can simply perform gradient-based approximate minimization (such as minibatch gradient descent) on the negative log-likelihood  $-\log p_{\text{decoder}}(\mathbf{x} \mid \mathbf{h})$ . So long as the encoder is deterministic, the denoising autoencoder is a feedforward network and may be trained with exactly the same techniques as any other feedforward network.

We can therefore view the DAE as performing stochastic gradient descent on the following expectation:

$$-\mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}(\mathbf{x})} \mathbb{E}_{\tilde{\mathbf{x}} \sim C(\tilde{\mathbf{x}} \mid \mathbf{x})} \log p_{\text{decoder}}(\mathbf{x} \mid \mathbf{h} = f(\tilde{\mathbf{x}})) \quad (14.14)$$

where  $\hat{p}_{\text{data}}(\mathbf{x})$  is the training distribution.

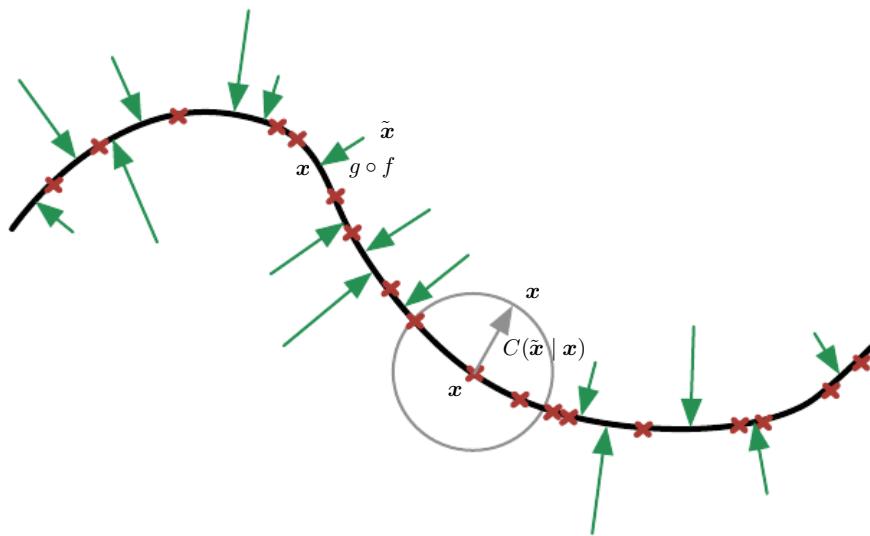


Figure 14.4: A denoising autoencoder is trained to map a corrupted data point  $\tilde{\mathbf{x}}$  back to the original data point  $\mathbf{x}$ . We illustrate training examples  $\mathbf{x}$  as red crosses lying near a low-dimensional manifold illustrated with the bold black line. We illustrate the corruption process  $C(\tilde{\mathbf{x}} | \mathbf{x})$  with a gray circle of equiprobable corruptions. A gray arrow demonstrates how one training example is transformed into one sample from this corruption process. When the denoising autoencoder is trained to minimize the average of squared errors  $\|g(f(\tilde{\mathbf{x}})) - \mathbf{x}\|^2$ , the reconstruction  $g(f(\tilde{\mathbf{x}}))$  estimates  $\mathbb{E}_{\mathbf{x}, \tilde{\mathbf{x}} \sim p_{\text{data}}(\mathbf{x})} C(\tilde{\mathbf{x}} | \mathbf{x})[\mathbf{x} | \tilde{\mathbf{x}}]$ . The vector  $g(f(\tilde{\mathbf{x}})) - \tilde{\mathbf{x}}$  points approximately towards the nearest point on the manifold, since  $g(f(\tilde{\mathbf{x}}))$  estimates the center of mass of the clean points  $\mathbf{x}$  which could have given rise to  $\tilde{\mathbf{x}}$ . The autoencoder thus learns a vector field  $g(f(\mathbf{x})) - \mathbf{x}$  indicated by the green arrows. This vector field estimates the score  $\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})$  up to a multiplicative factor that is the average root mean square reconstruction error.

### 14.5.1 Estimating the Score

Score matching (Hyvärinen, 2005) is an alternative to maximum likelihood. It provides a consistent estimator of probability distributions based on encouraging the model to have the same *score* as the data distribution at every training point  $\mathbf{x}$ . In this context, the score is a particular gradient field:

$$\nabla_{\mathbf{x}} \log p(\mathbf{x}). \quad (14.15)$$

Score matching is discussed further in Sec. 18.4. For the present discussion regarding autoencoders, it is sufficient to understand that learning the gradient field of  $\log p_{\text{data}}$  is one way to learn the structure of  $p_{\text{data}}$  itself.

A very important property of DAEs is that their training criterion (with conditionally Gaussian  $p(\mathbf{x} | \mathbf{h})$ ) makes the autoencoder learn a vector field ( $g(f(\mathbf{x})) - \mathbf{x}$ ) that estimates the score of the data distribution. This is illustrated in Fig. 14.4.

Denoising training of a specific kind of autoencoder (sigmoidal hidden units, linear reconstruction units) using Gaussian noise and mean squared error as the reconstruction cost is equivalent (Vincent, 2011) to training a specific kind of undirected probabilistic model called an RBM with Gaussian visible units. This kind of model will be described in detail in Sec. 20.5.1; for the present discussion it suffices to know that it is a model that provides an explicit  $p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$ . When the RBM is trained using *denoising score matching* (Kingma and LeCun, 2010), its learning algorithm is equivalent to denoising training in the corresponding autoencoder. With a fixed noise level, regularized score matching is not a consistent estimator; it instead recovers a blurred version of the distribution. However, if the noise level is chosen to approach 0 when the number of examples approaches infinity, then consistency is recovered. Denoising score matching is discussed in more detail in Sec. 18.5.

Other connections between autoencoders and RBMs exist. Score matching applied to RBMs yields a cost function that is identical to reconstruction error combined with a regularization term similar to the contractive penalty of the CAE (Swersky *et al.*, 2011). Bengio and Delalleau (2009) showed that an autoencoder gradient provides an approximation to contrastive divergence training of RBMs.

For continuous-valued  $\mathbf{x}$ , the denoising criterion with Gaussian corruption and reconstruction distribution yields an estimator of the score that is applicable to general encoder and decoder parametrizations (Alain and Bengio, 2013). This means a generic encoder-decoder architecture may be made to estimate the score

by training with the squared error criterion

$$\|g(f(\tilde{\mathbf{x}})) - \mathbf{x}\|^2 \quad (14.16)$$

and corruption

$$C(\tilde{\mathbf{x}} = \tilde{\mathbf{x}}|\mathbf{x}) = \mathcal{N}(\tilde{\mathbf{x}}; \mu = \mathbf{x}, \Sigma = \sigma^2 I) \quad (14.17)$$

with noise variance  $\sigma^2$ . See Fig. 14.5 for an illustration of how this works.

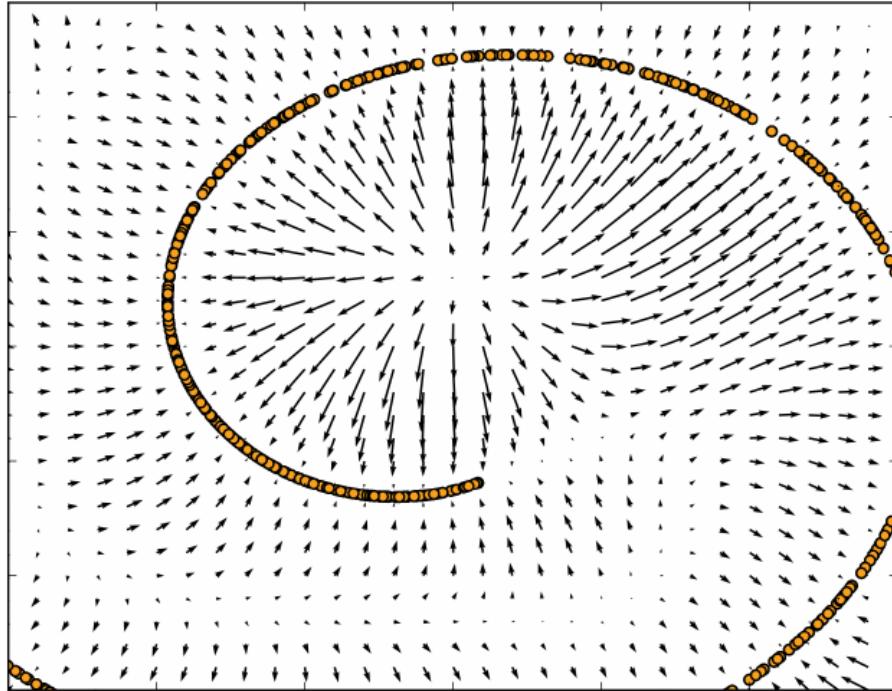


Figure 14.5: Vector field learned by a denoising autoencoder around a 1-D curved manifold near which the data concentrates in a 2-D space. Each arrow is proportional to the reconstruction minus input vector of the autoencoder and points towards higher probability according to the implicitly estimated probability distribution. The vector field has zeros at both maxima of the estimated density function (on the data manifolds) and at minima of that density function. For example, the spiral arm forms a one-dimensional manifold of local maxima that are connected to each other. Local minima appear near the middle of the gap between two arms. When the norm of reconstruction error (shown by the length of the arrows) is large, it means that probability can be significantly increased by moving in the direction of the arrow, and that is mostly the case in places of low probability. The autoencoder maps these low probability points to higher probability reconstructions. Where probability is maximal, the arrows shrink because the reconstruction becomes more accurate.

In general, there is no guarantee that the reconstruction  $g(f(\mathbf{x}))$  minus the input  $\mathbf{x}$  corresponds to the gradient of any function, let alone to the score. That is

why the early results (Vincent, 2011) are specialized to particular parametrizations where  $g(f(\mathbf{x})) - \mathbf{x}$  may be obtained by taking the derivative of another function. Kamyshanska and Memisevic (2015) generalized the results of Vincent (2011) by identifying a family of shallow autoencoders such that  $g(f(\mathbf{x})) - \mathbf{x}$  corresponds to a score for all members of the family.

So far we have described only how the denoising autoencoder learns to represent a probability distribution. More generally, one may want to use the autoencoder as a generative model and draw samples from this distribution. This will be described later, in Sec. 20.11.

#### 14.5.1.1 Historical Perspective

The idea of using MLPs for denoising dates back to the work of LeCun (1987) and Gallinari *et al.* (1987). Behnke (2001) also used recurrent networks to denoise images. Denoising autoencoders are, in some sense, just MLPs trained to denoise. However, the name “denoising autoencoder” refers to a model that is intended not merely to learn to denoise its input but to learn a good internal representation as a side effect of learning to denoise. This idea came much later (Vincent *et al.*, 2008, 2010). The learned representation may then be used to pretrain a deeper unsupervised network or a supervised network. Like sparse autoencoders, sparse coding, contractive autoencoders and other regularized autoencoders, the motivation for DAEs was to allow the learning of a very high-capacity encoder while preventing the encoder and decoder from learning a useless identity function.

Prior to the introduction of the modern DAE, Inayoshi and Kurita (2005) explored some of the same goals with some of the same methods. Their approach minimizes reconstruction error in addition to a supervised objective while injecting noise in the hidden layer of a supervised MLP, with the objective to improve generalization by introducing the reconstruction error and the injected noise. However, their method was based on a linear encoder and could not learn function families as powerful as can the modern DAE.

## 14.6 Learning Manifolds with Autoencoders

Like many other machine learning algorithms, autoencoders exploit the idea that data concentrates around a low-dimensional manifold or a small set of such manifolds, as described in Sec. 5.11.3. Some machine learning algorithms exploit this idea only insofar as that they learn a function that behaves correctly on the manifold but may have unusual behavior if given an input that is off the manifold.

Autoencoders take this idea further and aim to learn the structure of the manifold.

To understand how autoencoders do this, we must present some important characteristics of manifolds.

An important characterization of a manifold is the set of its *tangent planes*. At a point  $\mathbf{x}$  on a  $d$ -dimensional manifold, the tangent plane is given by  $d$  basis vectors that span the local directions of variation allowed on the manifold. As illustrated in Fig. 14.6, these local directions specify how one can change  $\mathbf{x}$  infinitesimally while staying on the manifold.

All autoencoder training procedures involve a compromise between two forces:

1. Learning a representation  $\mathbf{h}$  of a training example  $\mathbf{x}$  such that  $\mathbf{x}$  can be approximately recovered from  $\mathbf{h}$  through a decoder. The fact that  $\mathbf{x}$  is drawn from the training data is crucial, because it means the autoencoder need not successfully reconstruct inputs that are not probable under the data generating distribution.
2. Satisfying the constraint or regularization penalty. This can be an architectural constraint that limits the capacity of the autoencoder, or it can be a regularization term added to the reconstruction cost. These techniques generally prefer solutions that are less sensitive to the input.

Clearly, neither force alone would be useful—copying the input to the output is not useful on its own, nor is ignoring the input. Instead, the two forces together are useful because they force the hidden representation to capture information about the structure of the data generating distribution. The important principle is that the autoencoder can afford to represent *only the variations that are needed to reconstruct training examples*. If the data generating distribution concentrates near a low-dimensional manifold, this yields representations that implicitly capture a local coordinate system for this manifold: only the variations tangent to the manifold around  $\mathbf{x}$  need to correspond to changes in  $\mathbf{h} = f(\mathbf{x})$ . Hence the encoder learns a mapping from the input space  $\mathbf{x}$  to a representation space, a mapping that is only sensitive to changes along the manifold directions, but that is insensitive to changes orthogonal to the manifold.

A one-dimensional example is illustrated in Fig. 14.7, showing that by making the reconstruction function insensitive to perturbations of the input around the data points we recover the manifold structure.

To understand why autoencoders are useful for manifold learning, it is instructive to compare them to other approaches. What is most commonly learned to characterize a manifold is a *representation* of the data points on (or near) the

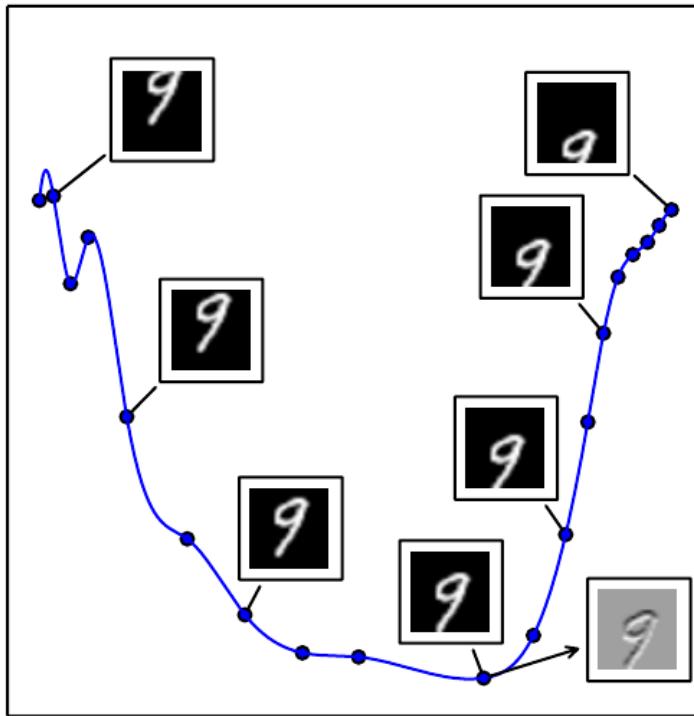


Figure 14.6: An illustration of the concept of a tangent hyperplane. Here we create a one-dimensional manifold in 784-dimensional space. We take an MNIST image with 784 pixels and transform it by translating it vertically. The amount of vertical translation defines a coordinate along a one-dimensional manifold that traces out a curved path through image space. This plot shows a few points along this manifold. For visualization, we have projected the manifold into two dimensional space using PCA. An  $n$ -dimensional manifold has an  $n$ -dimensional tangent plane at every point. This tangent plane touches the manifold exactly at that point and is oriented parallel to the surface at that point. It defines the space of directions in which it is possible to move while remaining on the manifold. This one-dimensional manifold has a single tangent line. We indicate an example tangent line at one point, with an image showing how this tangent direction appears in image space. Gray pixels indicate pixels that do not change as we move along the tangent line, white pixels indicate pixels that brighten, and black pixels indicate pixels that darken.

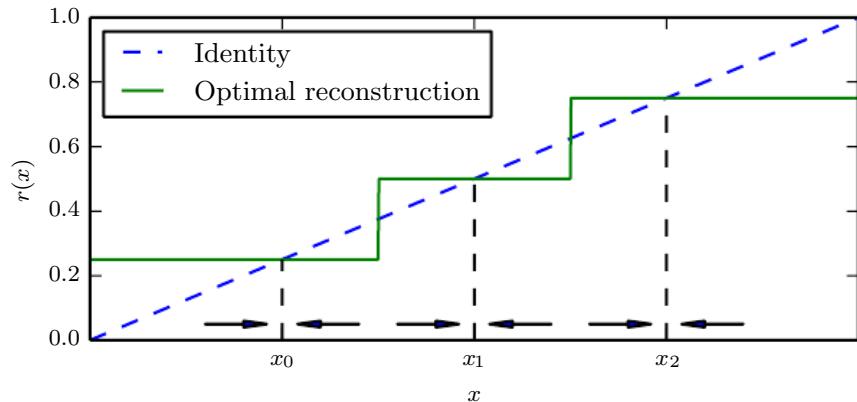


Figure 14.7: If the autoencoder learns a reconstruction function that is invariant to small perturbations near the data points, it captures the manifold structure of the data. Here the manifold structure is a collection of 0-dimensional manifolds. The dashed diagonal line indicates the identity function target for reconstruction. The optimal reconstruction function crosses the identity function wherever there is a data point. The horizontal arrows at the bottom of the plot indicate the  $r(\mathbf{x}) - \mathbf{x}$  reconstruction direction vector at the base of the arrow, in input space, always pointing towards the nearest “manifold” (a single datapoint, in the 1-D case). The denoising autoencoder explicitly tries to make the derivative of the reconstruction function  $r(\mathbf{x})$  small around the data points. The contractive autoencoder does the same for the encoder. Although the derivative of  $r(\mathbf{x})$  is asked to be small around the data points, it can be large between the data points. The space between the data points corresponds to the region between the manifolds, where the reconstruction function must have a large derivative in order to map corrupted points back onto the manifold.

manifold. Such a representation for a particular example is also called its embedding. It is typically given by a low-dimensional vector, with less dimensions than the “ambient” space of which the manifold is a low-dimensional subset. Some algorithms (non-parametric manifold learning algorithms, discussed below) directly learn an embedding for each training example, while others learn a more general mapping, sometimes called an encoder, or representation function, that maps any point in the ambient space (the input space) to its embedding.

Manifold learning has mostly focused on unsupervised learning procedures that attempt to capture these manifolds. Most of the initial machine learning research on learning nonlinear manifolds has focused on *non-parametric* methods based on the *nearest-neighbor graph*. This graph has one node per training example and edges connecting near neighbors to each other. These methods (Schölkopf *et al.*, 1998; Roweis and Saul, 2000; Tenenbaum *et al.*, 2000; Brand, 2003; Belkin and

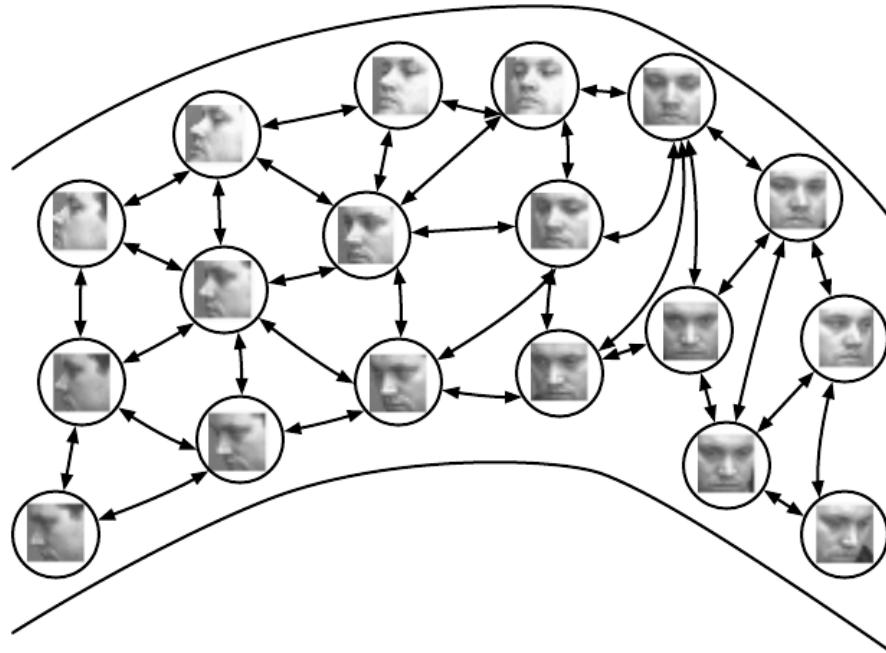


Figure 14.8: Non-parametric manifold learning procedures build a nearest neighbor graph whose nodes are training examples and arcs connect nearest neighbors. Various procedures can thus obtain the tangent plane associated with a neighborhood of the graph as well as a coordinate system that associates each training example with a real-valued vector position, or *embedding*. It is possible to generalize such a representation to new examples by a form of interpolation. So long as the number of examples is large enough to cover the curvature and twists of the manifold, these approaches work well. Images from the QMUL Multiview Face Dataset (Gong *et al.*, 2000).

Niyogi, 2003; Donoho and Grimes, 2003; Weinberger and Saul, 2004; Hinton and Roweis, 2003; van der Maaten and Hinton, 2008) associate each of nodes with a tangent plane that spans the directions of variations associated with the difference vectors between the example and its neighbors, as illustrated in Fig. 14.8.

A global coordinate system can then be obtained through an optimization or solving a linear system. Fig. 14.9 illustrates how a manifold can be tiled by a large number of locally linear Gaussian-like patches (or “pancakes,” because the Gaussians are flat in the tangent directions).

However, there is a fundamental difficulty with such local non-parametric approaches to manifold learning, raised in Bengio and Monperrus (2005): if the manifolds are not very smooth (they have many peaks and troughs and twists), one may need a very large number of training examples to cover each one of these variations, with no chance to generalize to unseen variations. Indeed, these methods

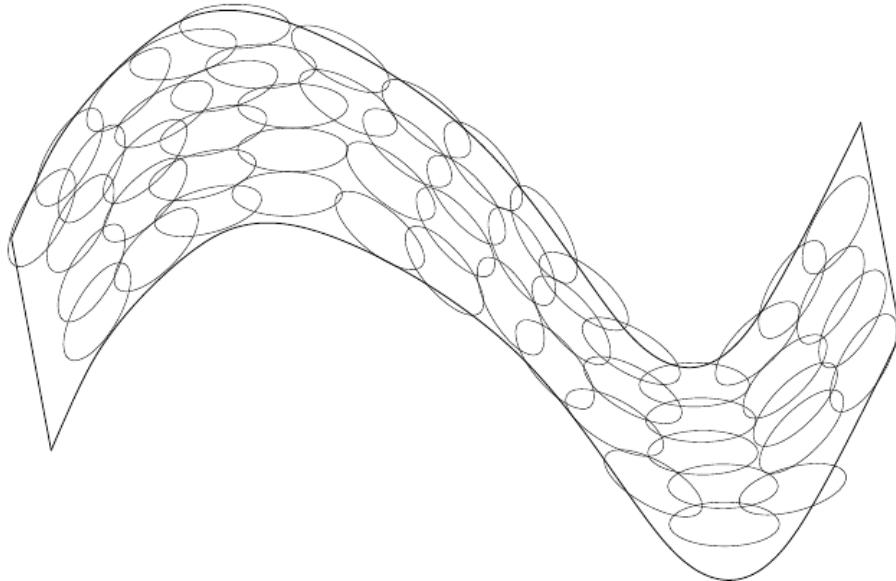


Figure 14.9: If the tangent planes (see Fig. 14.6) at each location are known, then they can be tiled to form a global coordinate system or a density function. Each local patch can be thought of as a local Euclidean coordinate system or as a locally flat Gaussian, or “pancake”, with a very small variance in the directions orthogonal to the pancake and a very large variance in the directions defining the coordinate system on the pancake. A mixture of these Gaussians provides an estimated density function, as in the manifold Parzen window algorithm (Vincent and Bengio, 2003) or its non-local neural-net based variant (Bengio *et al.*, 2006c).

can only generalize the shape of the manifold by interpolating between neighboring examples. Unfortunately, the manifolds involved in AI problems can have very complicated structure that can be difficult to capture from only local interpolation. Consider for example the manifold resulting from translation shown in Fig. 14.6. If we watch just one coordinate within the input vector,  $x_i$ , as the image is translated, we will observe that one coordinate encounters a peak or a trough in its value once for every peak or trough in brightness in the image. In other words, the complexity of the patterns of brightness in an underlying image template drives the complexity of the manifolds that are generated by performing simple image transformations. This motivates the use of distributed representations and deep learning for capturing manifold structure.

## 14.7 Contractive Autoencoders

The contractive autoencoder (Rifai *et al.*, 2011a,b) introduces an explicit regularizer on the code  $\mathbf{h} = f(\mathbf{x})$ , encouraging the derivatives of  $f$  to be as small as possible:

$$\Omega(\mathbf{h}) = \lambda \left\| \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right\|_F^2. \quad (14.18)$$

The penalty  $\Omega(\mathbf{h})$  is the squared Frobenius norm (sum of squared elements) of the Jacobian matrix of partial derivatives associated with the encoder function.

There is a connection between the denoising autoencoder and the contractive autoencoder: Alain and Bengio (2013) showed that in the limit of small Gaussian input noise, the denoising reconstruction error is equivalent to a contractive penalty on the reconstruction function that maps  $\mathbf{x}$  to  $\mathbf{r} = g(f(\mathbf{x}))$ . In other words, denoising autoencoders make the reconstruction function resist small but finite-sized perturbations of the input, while contractive autoencoders make the feature extraction function resist infinitesimal perturbations of the input. When using the Jacobian-based contractive penalty to pretrain features  $f(\mathbf{x})$  for use with a classifier, the best classification accuracy usually results from applying the contractive penalty to  $f(\mathbf{x})$  rather than to  $g(f(\mathbf{x}))$ . A contractive penalty on  $f(\mathbf{x})$  also has close connections to score matching, as discussed in Sec. 14.5.1.

The name *contractive* arises from the way that the CAE warps space. Specifically, because the CAE is trained to resist perturbations of its input, it is encouraged to map a neighborhood of input points to a smaller neighborhood of output points. We can think of this as contracting the input neighborhood to a smaller output neighborhood.

To clarify, the CAE is contractive only locally—all perturbations of a training point  $\mathbf{x}$  are mapped near to  $f(\mathbf{x})$ . Globally, two different points  $\mathbf{x}$  and  $\mathbf{x}'$  may be mapped to  $f(\mathbf{x})$  and  $f(\mathbf{x}')$  points that are farther apart than the original points. It is plausible that  $f$  be expanding in-between or far from the data manifolds (see for example what happens in the 1-D toy example of Fig. 14.7). When the  $\Omega(\mathbf{h})$  penalty is applied to sigmoidal units, one easy way to shrink the Jacobian is to make the sigmoid units saturate to 0 or 1. This encourages the CAE to encode input points with extreme values of the sigmoid that may be interpreted as a binary code. It also ensures that the CAE will spread its code values throughout most of the hypercube that its sigmoidal hidden units can span.

We can think of the Jacobian matrix  $\mathbf{J}$  at a point  $\mathbf{x}$  as approximating the nonlinear encoder  $f(\mathbf{x})$  as being a linear operator. This allows us to use the word “contractive” more formally. In the theory of linear operators, a linear operator

is said to be contractive if the norm of  $\mathbf{J}\mathbf{x}$  remains less than or equal to 1 for all unit-norm  $\mathbf{x}$ . In other words,  $\mathbf{J}$  is contractive if it shrinks the unit sphere. We can think of the CAE as penalizing the Frobenius norm of the local linear approximation of  $f(\mathbf{x})$  at every training point  $\mathbf{x}$  in order to encourage each of these local linear operator to become a contraction.

As described in Sec. 14.6, regularized autoencoders learn manifolds by balancing two opposing forces. In the case of the CAE, these two forces are reconstruction error and the contractive penalty  $\Omega(\mathbf{h})$ . Reconstruction error alone would encourage the CAE to learn an identity function. The contractive penalty alone would encourage the CAE to learn features that are constant with respect to  $\mathbf{x}$ . The compromise between these two forces yields an autoencoder whose derivatives  $\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}$  are mostly tiny. Only a small number of hidden units, corresponding to a small number of directions in the input, may have significant derivatives.

The goal of the CAE is to learn the manifold structure of the data. Directions  $\mathbf{x}$  with large  $\mathbf{J}\mathbf{x}$  rapidly change  $\mathbf{h}$ , so these are likely to be directions which approximate the tangent planes of the manifold. Experiments by Rifai *et al.* (2011a) and Rifai *et al.* (2011b) show that training the CAE results in most singular values of  $\mathbf{J}$  dropping below 1 in magnitude and therefore becoming contractive. However, some singular values remain above 1, because the reconstruction error penalty encourages the CAE to encode the directions with the most local variance. The directions corresponding to the largest singular values are interpreted as the tangent directions that the contractive autoencoder has learned. Ideally, these tangent directions should correspond to real variations in the data. For example, a CAE applied to images should learn tangent vectors that show how the image changes as objects in the image gradually change pose, as shown in Fig. 14.6. Visualizations of the experimentally obtained singular vectors do seem to correspond to meaningful transformations of the input image, as shown in Fig. 14.10.

One practical issue with the CAE regularization criterion is that although it is cheap to compute in the case of a single hidden layer autoencoder, it becomes much more expensive in the case of deeper autoencoders. The strategy followed by Rifai *et al.* (2011a) is to separately train a series of single-layer autoencoders, each trained to reconstruct the previous autoencoder's hidden layer. The composition of these autoencoders then forms a deep autoencoder. Because each layer was separately trained to be locally contractive, the deep autoencoder is contractive as well. The result is not the same as what would be obtained by jointly training the entire architecture with a penalty on the Jacobian of the deep model, but it captures many of the desirable qualitative characteristics.

Another practical issue is that the contraction penalty can obtain useless results

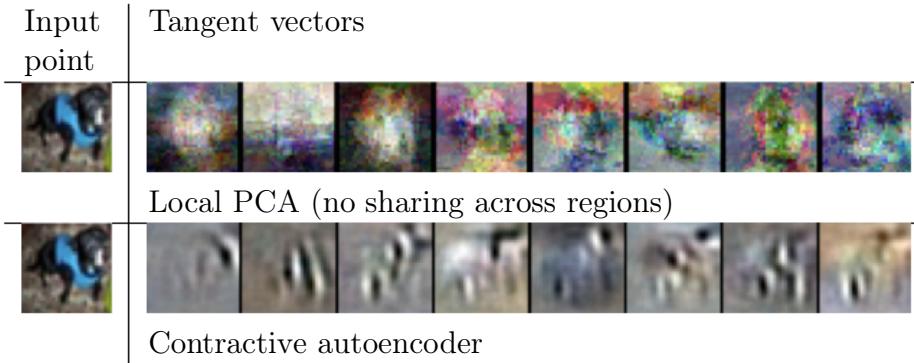


Figure 14.10: Illustration of tangent vectors of the manifold estimated by local PCA and by a contractive autoencoder. The location on the manifold is defined by the input image of a dog drawn from the CIFAR-10 dataset. The tangent vectors are estimated by the leading singular vectors of the Jacobian matrix  $\frac{\partial \mathbf{h}}{\partial \mathbf{x}}$  of the input-to-code mapping. Although both local PCA and the CAE can capture local tangents, the CAE is able to form more accurate estimates from limited training data because it exploits parameter sharing across different locations that share a subset of active hidden units. The CAE tangent directions typically correspond to moving or changing parts of the object (such as the head or legs).

if we do not impose some sort of scale on the decoder. For example, the encoder could consist of multiplying the input by a small constant  $\epsilon$  and the decoder could consist of dividing the code by  $\epsilon$ . As  $\epsilon$  approaches 0, the encoder drives the contractive penalty  $\Omega(\mathbf{h})$  to approach 0 without having learned anything about the distribution. Meanwhile, the decoder maintains perfect reconstruction. In [Rifai et al. \(2011a\)](#), this is prevented by tying the weights of  $f$  and  $g$ . Both  $f$  and  $g$  are standard neural network layers consisting of an affine transformation followed by an element-wise nonlinearity, so it is straightforward to set the weight matrix of  $g$  to be the transpose of the weight matrix of  $f$ .

## 14.8 Predictive Sparse Decomposition

*Predictive sparse decomposition* (PSD) is a model that is a hybrid of sparse coding and parametric autoencoders ([Kavukcuoglu et al., 2008](#)). A parametric encoder is trained to predict the output of iterative inference. PSD has been applied to unsupervised feature learning for object recognition in images and video ([Kavukcuoglu et al., 2009, 2010; Jarrett et al., 2009; Farabet et al., 2011](#)), as well as for audio ([Henaff et al., 2011](#)). The model consists of an encoder  $f(\mathbf{x})$  and a decoder  $g(\mathbf{h})$  that are both parametric. During training,  $\mathbf{h}$  is controlled by the

optimization algorithm. Training proceeds by minimizing

$$\|\mathbf{x} - g(\mathbf{h})\|^2 + \lambda|\mathbf{h}|_1 + \gamma\|\mathbf{h} - f(\mathbf{x})\|^2. \quad (14.19)$$

Like in sparse coding, the training algorithm alternates between minimization with respect to  $\mathbf{h}$  and minimization with respect to the model parameters. Minimization with respect to  $\mathbf{h}$  is fast because  $f(\mathbf{x})$  provides a good initial value of  $\mathbf{h}$  and the cost function constrains  $\mathbf{h}$  to remain near  $f(\mathbf{x})$  anyway. Simple gradient descent can obtain reasonable values of  $\mathbf{h}$  in as few as ten steps.

The training procedure used by PSD is different from first training a sparse coding model and then training  $f(\mathbf{x})$  to predict the values of the sparse coding features. The PSD training procedure regularizes the decoder to use parameters for which  $f(\mathbf{x})$  can infer good code values.

Predictive sparse coding is an example of *learned approximate inference*. In Sec. 19.5, this topic is developed further. The tools presented in Chapter 19 make it clear that PSD can be interpreted as training a directed sparse coding probabilistic model by maximizing a lower bound on the log-likelihood of the model.

In practical applications of PSD, the iterative optimization is only used during training. The parametric encoder  $f$  is used to compute the learned features when the model is deployed. Evaluating  $f$  is computationally inexpensive compared to inferring  $\mathbf{h}$  via gradient descent. Because  $f$  is a differentiable parametric function, PSD models may be stacked and used to initialize a deep network to be trained with another criterion.

## 14.9 Applications of Autoencoders

Autoencoders have been successfully applied to dimensionality reduction and information retrieval tasks. Dimensionality reduction was one of the first applications of representation learning and deep learning. It was one of the early motivations for studying autoencoders. For example, Hinton and Salakhutdinov (2006) trained a stack of RBMs and then used their weights to initialize a deep autoencoder with gradually smaller hidden layers, culminating in a bottleneck of 30 units. The resulting code yielded less reconstruction error than PCA into 30 dimensions and the learned representation was qualitatively easier to interpret and relate to the underlying categories, with these categories manifesting as well-separated clusters.

Lower-dimensional representations can improve performance on many tasks, such as classification. Models of smaller spaces consume less memory and runtime. Many forms of dimensionality reduction place semantically related examples near

each other, as observed by Salakhutdinov and Hinton (2007b) and Torralba *et al.* (2008). The hints provided by the mapping to the lower-dimensional space aid generalization.

One task that benefits even more than usual from dimensionality reduction is *information retrieval*, the task of finding entries in a database that resemble a query entry. This task derives the usual benefits from dimensionality reduction that other tasks do, but also derives the additional benefit that search can become extremely efficient in certain kinds of low dimensional spaces. Specifically, if we train the dimensionality reduction algorithm to produce a code that is low-dimensional and **binary**, then we can store all database entries in a hash table mapping binary code vectors to entries. This hash table allows us to perform information retrieval by returning all database entries that have the same binary code as the query. We can also search over slightly less similar entries very efficiently, just by flipping individual bits from the encoding of the query. This approach to information retrieval via dimensionality reduction and binarization is called *semantic hashing* (Salakhutdinov and Hinton, 2007b, 2009b), and has been applied to both textual input (Salakhutdinov and Hinton, 2007b, 2009b) and images (Torralba *et al.*, 2008; Weiss *et al.*, 2008; Krizhevsky and Hinton, 2011).

To produce binary codes for semantic hashing, one typically uses an encoding function with sigmoids on the final layer. The sigmoid units must be trained to be saturated to nearly 0 or nearly 1 for all input values. One trick that can accomplish this is simply to inject additive noise just before the sigmoid nonlinearity during training. The magnitude of the noise should increase over time. To fight that noise and preserve as much information as possible, the network must increase the magnitude of the inputs to the sigmoid function, until saturation occurs.

The idea of learning a hashing function has been further explored in several directions, including the idea of training the representations so as to optimize a loss more directly linked to the task of finding nearby examples in the hash table (Norouzi and Fleet, 2011).

## Chapter 15

# Representation Learning

In this chapter, we first discuss what it means to learn representations and how the notion of representation can be useful to design deep architectures. We discuss how learning algorithms share statistical strength across different tasks, including using information from unsupervised tasks to perform supervised tasks. Shared representations are useful to handle multiple modalities or domains, or to transfer learned knowledge to tasks for which few or no examples are given but a task representation exists. Finally, we step back and argue about the reasons for the success of representation learning, starting with the theoretical advantages of distributed representations (Hinton *et al.*, 1986) and deep representations and ending with the more general idea of underlying assumptions about the data generating process, in particular about underlying causes of the observed data.

Many information processing tasks can be very easy or very difficult depending on how the information is represented. This is a general principle applicable to daily life, computer science in general, and to machine learning. For example, it is straightforward for a person to divide 210 by 6 using long division. The task becomes considerably less straightforward if it is instead posed using the Roman numeral representation of the numbers. Most modern people asked to divide CCX by VI would begin by converting the numbers to the Arabic numeral representation, permitting long division procedures that make use of the place value system. More concretely, we can quantify the asymptotic runtime of various operations using appropriate or inappropriate representations. For example, inserting a number into the correct position in a sorted list of numbers is an  $O(n)$  operation if the list is represented as a linked list, but only  $O(\log n)$  if the list is represented as a red-black tree.

In the context of machine learning, what makes one representation better than

another? Generally speaking, a good representation is one that makes a subsequent learning task easier. The choice of representation will usually depend on the choice of the subsequent learning task.

We can think of feedforward networks trained by supervised learning as performing a kind of representation learning. Specifically, the last layer of the network is typically a linear classifier, such as a softmax regression classifier. The rest of the network learns to provide a representation to this classifier. Training with a supervised criterion naturally leads to the representation at every hidden layer (but more so near the top hidden layer) taking on properties that make the classification task easier. For example, classes that were not linearly separable in the input features may become linearly separable in the last hidden layer. In principle, the last layer could be another kind of model, such as a nearest neighbor classifier ([Salakhutdinov and Hinton, 2007a](#)). The features in the penultimate layer should learn different properties depending on the type of the last layer.

Supervised training of feedforward networks does not involve explicitly imposing any condition on the learned intermediate features. Other kinds of representation learning algorithms are often explicitly designed to shape the representation in some particular way. For example, suppose we want to learn a representation that makes density estimation easier. Distributions with more independences are easier to model, so we could design an objective function that encourages the elements of the representation vector  $\mathbf{h}$  to be independent. Just like supervised networks, unsupervised deep learning algorithms have a main training objective but also learn a representation as a side effect. Regardless of how a representation was obtained, it can be used for another task. Alternatively, multiple tasks (some supervised, some unsupervised) can be learned together with some shared internal representation.

Most representation learning problems face a tradeoff between preserving as much information about the input as possible and attaining nice properties (such as independence).

Representation learning is particularly interesting because it provides one way to perform unsupervised and semi-supervised learning. We often have very large amounts of unlabeled training data and relatively little labeled training data. Training with supervised learning techniques on the labeled subset often results in severe overfitting. Semi-supervised learning offers the chance to resolve this overfitting problem by also learning from the unlabeled data. Specifically, we can learn good representations for the unlabeled data, and then use these representations to solve the supervised learning task.

Humans and animals are able to learn from very few labeled examples. We do

not yet know how this is possible. Many factors could explain improved human performance—for example, the brain may use very large ensembles of classifiers or Bayesian inference techniques. One popular hypothesis is that the brain is able to leverage unsupervised or semi-supervised learning. There are many ways to leverage unlabeled data. In this chapter, we focus on the hypothesis that the unlabeled data can be used to learn a good representation.

## 15.1 Greedy Layer-Wise Unsupervised Pretraining

Unsupervised learning played a key historical role in the revival of deep neural networks, allowing for the first time to train a deep supervised network without requiring architectural specializations like convolution or recurrence. We call this procedure *unsupervised pretraining*, or more precisely, *greedy layer-wise unsupervised pretraining*. This procedure is a canonical example of how a representation learned for one task (unsupervised learning, trying to capture the shape of the input distribution) can sometimes be useful for another task (supervised learning with the same input domain).

Greedy layer-wise unsupervised pretraining relies on a single-layer representation learning algorithm such as an RBM, a single-layer autoencoder, a sparse coding model, or another model that learns latent representations. Each layer is pretrained using unsupervised learning, taking the output of the previous layer and producing as output a new representation of the data, whose distribution (or its relation to other variables such as categories to predict) is hopefully simpler. See Algorithm 15.1 for a formal description.

Greedy layer-wise training procedures based on unsupervised criteria have long been used to sidestep the difficulty of jointly training the layers of a deep neural net for a supervised task. This approach dates back at least as far as the Neocognitron (Fukushima, 1975). The deep learning renaissance of 2006 began with the discovery that this greedy learning procedure could be used to find a good initialization for a joint learning procedure over all the layers, and that this approach could be used to successfully train even fully connected architectures (Hinton *et al.*, 2006; Hinton and Salakhutdinov, 2006; Hinton, 2006; Bengio *et al.*, 2007; Ranzato *et al.*, 2007a). Prior to this discovery, only convolutional deep networks or networks whose depth resulted from recurrence were regarded as feasible to train. Today, we now know that greedy layer-wise pretraining is not required to train fully connected deep architectures, but the unsupervised pretraining approach was the first method to succeed.

Greedy layer-wise pretraining is called *greedy* because it is a *greedy algorithm*,

meaning that it optimizes each piece of the solution independently, one piece at a time, rather than jointly optimizing all pieces. It is called *layer-wise* because these independent pieces are the layers of the network. Specifically, greedy layer-wise pretraining proceeds one layer at a time, training the  $k$ -th layer while keeping the previous ones fixed. In particular, the lower layers (which are trained first) are not adapted after the upper layers are introduced. It is called *unsupervised* because each layer is trained with an unsupervised representation learning algorithm. However it is also called *pretraining*, because it is supposed to be only a first step before a joint training algorithm is applied to *fine-tune* all the layers together. In the context of a supervised learning task, it can be viewed as a regularizer (in some experiments, pretraining decreases test error without decreasing training error) and a form of parameter initialization.

It is common to use the word “pretraining” to refer not only to the pretraining stage itself but to the entire two phase protocol that combines the pretraining phase and a supervised learning phase. The supervised learning phase may involve training a simple classifier on top of the features learned in the pretraining phase, or it may involve supervised fine-tuning of the entire network learned in the pretraining phase. No matter what kind of unsupervised learning algorithm or what model type is employed, in the vast majority of cases, the overall training scheme is nearly the same. While the choice of unsupervised learning algorithm will obviously impact the details, most applications of unsupervised pretraining follow this basic protocol.

Greedy layer-wise unsupervised pretraining can also be used as initialization for other unsupervised learning algorithms, such as deep autoencoders ([Hinton and Salakhutdinov, 2006](#)) and probabilistic models with many layers of latent variables. Such models include deep belief networks ([Hinton et al., 2006](#)) and deep Boltzmann machines ([Salakhutdinov and Hinton, 2009a](#)). These deep generative models will be described in Chapter [20](#).

As discussed in Sec. [8.7.4](#), it is also possible to have greedy layer-wise *supervised* pretraining. This builds on the premise that training a shallow network is easier than training a deep one, which seems to have been validated in several contexts ([Erhan et al., 2010](#)).

### 15.1.1 When and Why Does Unsupervised Pretraining Work?

On many tasks, greedy layer-wise unsupervised pretraining can yield substantial improvements in test error for classification tasks. This observation was responsible for the renewed interest in deep neural networks starting in 2006 ([Hinton et al.](#),

**Algorithm 15.1 Greedy layer-wise unsupervised pretraining protocol.**

Given the following: Unsupervised feature learning algorithm  $\mathcal{L}$ , which takes a training set of examples and returns an encoder or feature function  $f$ . The raw input data is  $\mathbf{X}$ , with one row per example and  $f^{(1)}(\mathbf{X})$  is the output of the first stage encoder on  $\mathbf{X}$  and the dataset used by the second level unsupervised feature learner. In the case where fine-tuning is performed, we use a learner  $\mathcal{T}$  which takes an initial function  $f$ , input examples  $\mathbf{X}$  (and in the supervised fine-tuning case, associated targets  $\mathbf{Y}$ ), and returns a tuned function. The number of stages is  $m$ .

---

```
 $f \leftarrow$  Identity function  
 $\tilde{\mathbf{X}} = \mathbf{X}$   
for  $k = 1, \dots, m$  do  
     $f^{(k)} = \mathcal{L}(\tilde{\mathbf{X}})$   
     $f \leftarrow f^{(k)} \circ f$   
     $\tilde{\mathbf{X}} \leftarrow f^{(k)}(\tilde{\mathbf{X}})$   
end for  
if fine-tuning then  
     $f \leftarrow \mathcal{T}(f, \mathbf{X}, \mathbf{Y})$   
end if  
Return  $f$ 
```

---

2006; Bengio *et al.*, 2007; Ranzato *et al.*, 2007a). On many other tasks, however, unsupervised pretraining either does not confer a benefit or even causes noticeable harm. Ma *et al.* (2015) studied the effect of pretraining on machine learning models for chemical activity prediction and found that, on average, pretraining was slightly harmful, but for many tasks was significantly helpful. Because unsupervised pretraining is sometimes helpful but often harmful it is important to understand when and why it works in order to determine whether it is applicable to a particular task.

At the outset, it is important to clarify that most of this discussion is restricted to greedy unsupervised pretraining in particular. There are other, completely different paradigms for performing semi-supervised learning with neural networks, such as virtual adversarial training described in Sec. 7.13. It is also possible to train an autoencoder or generative model at the same time as the supervised model. Examples of this single-stage approach include the discriminative RBM (Larochelle and Bengio, 2008) and the ladder network (Rasmus *et al.*, 2015), in which the total objective is an explicit sum of the two terms (one using the labels and one only using the input).

Unsupervised pretraining combines two different ideas. First, it makes use of

the idea that the choice of initial parameters for a deep neural network can have a significant regularizing effect on the model (and, to a lesser extent, that it can improve optimization). Second, it makes use of the more general idea that learning about the input distribution can help to learn about the mapping from inputs to outputs.

Both of these ideas involve many complicated interactions between several parts of the machine learning algorithm that are not entirely understood.

The first idea, that the choice of initial parameters for a deep neural network can have a strong regularizing effect on its performance, is the least well understood. At the time that pretraining became popular, it was understood as initializing the model in a location that would cause it to approach one local minimum rather than another. Today, local minima are no longer considered to be a serious problem for neural network optimization. We now know that our standard neural network training procedures usually do not arrive at a critical point of any kind. It remains possible that pretraining initializes the model in a location that would otherwise be inaccessible—for example, a region that is surrounded by areas where the cost function varies so much from one example to another that minibatches give only a very noisy estimate of the gradient, or a region surrounded by areas where the Hessian matrix is so poorly conditioned that gradient descent methods must use very small steps. However, our ability to characterize exactly what aspects of the pretrained parameters are retained during the supervised training stage is limited. This is one reason that modern approaches typically use simultaneous unsupervised learning and supervised learning rather than two sequential stages. One may also avoid struggling with these complicated ideas about how optimization in the supervised learning stage preserves information from the unsupervised learning stage by simply freezing the parameters for the feature extractors and using supervised learning only to add a classifier on top of the learned features.

The other idea, that a learning algorithm can use information learned in the unsupervised phase to perform better in the supervised learning stage, is better understood. The basic idea is that some features that are useful for the unsupervised task may also be useful for the supervised learning task. For example, if we train a generative model of images of cars and motorcycles, it will need to know about wheels, and about how many wheels should be in an image. If we are fortunate, the representation of the wheels will take on a form that is easy for the supervised learner to access. This is not yet understood at a mathematical, theoretical level, so it is not always possible to predict which tasks will benefit from unsupervised learning in this way. Many aspects of this approach are highly dependent on the specific models used. For example, if we wish to add a linear classifier on

top of pretrained features, the features must make the underlying classes linearly separable. These properties often occur naturally but do not always do so. This is another reason that simultaneous supervised and unsupervised learning can be preferable—the constraints imposed by the output layer are naturally included from the start.

From the point of view of unsupervised pretraining as learning a representation, we can expect unsupervised pretraining to be more effective when the initial representation is poor. One key example of this is the use of word embeddings. Words represented by one-hot vectors are not very informative because every two distinct one-hot vectors are the same distance from each other (squared  $L^2$  distance of 2). Learned word embeddings naturally encode similarity between words by their distance from each other. Because of this, unsupervised pretraining is especially useful when processing words. It is less useful when processing images, perhaps because images already lie in a rich vector space where distances provide a low quality similarity metric.

From the point of view of unsupervised pretraining as a regularizer, we can expect unsupervised pretraining to be most helpful when the number of labeled examples is very small. Because the source of information added by unsupervised pretraining is the unlabeled data, we may also expect unsupervised pretraining to perform best when the number of unlabeled examples is very large. The advantage of semi-supervised learning via unsupervised pretraining with many unlabeled examples and few labeled examples was made particularly clear in 2011 with unsupervised pretraining winning two international transfer learning competitions ([Mesnil et al., 2011](#); [Goodfellow et al., 2011](#)), in settings where the number of labeled examples in the target task was small (from a handful to dozens of examples per class). These effects were also documented in carefully controlled experiments by [Paine et al. \(2014\)](#).

Other factors are likely to be involved. For example, unsupervised pretraining is likely to be most useful when the function to be learned is extremely complicated. Unsupervised learning differs from regularizers like weight decay because it does not bias the learner toward discovering a simple function but rather toward discovering feature functions that are useful for the unsupervised learning task. If the true underlying functions are complicated and shaped by regularities of the input distribution, unsupervised learning can be a more appropriate regularizer.

These caveats aside, we now analyze some success cases where unsupervised pretraining is known to cause an improvement, and explain what is known about why this improvement occurs. Unsupervised pretraining has usually been used to improve classifiers, and is usually most interesting from the point of view of

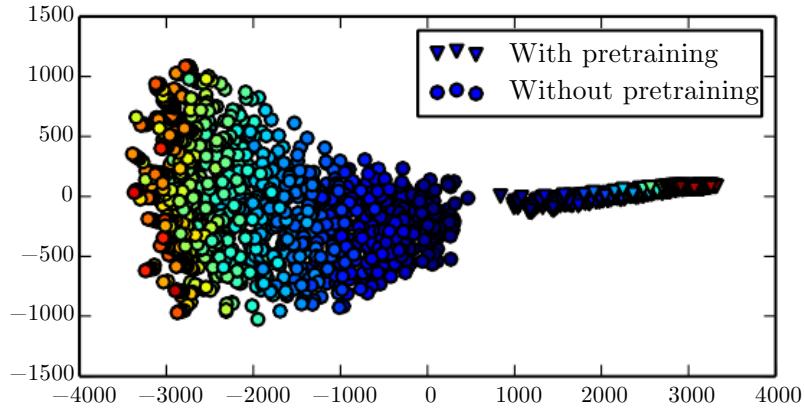


Figure 15.1: Visualization via nonlinear projection of the learning trajectories of different neural networks in **function space** (not parameter space, to avoid the issue of many-to-one mappings from parameter vectors to functions), with different random initializations and with or without unsupervised pretraining. Each point corresponds to a different neural network, at a particular time during its training process. This figure is adapted with permission from Erhan *et al.* (2010). A coordinate in function space is an infinite-dimensional vector associating every input  $\mathbf{x}$  with an output  $\mathbf{y}$ . Erhan *et al.* (2010) made a linear projection to high-dimensional space by concatenating the  $\mathbf{y}$  for many specific  $\mathbf{x}$  points. They then made a further nonlinear projection to 2-D by Isomap (Tenenbaum *et al.*, 2000). Color indicates time. All networks are initialized near the center of the plot (corresponding to the region of functions that produce approximately uniform distributions over the class  $y$  for most inputs). Over time, learning moves the function outward, to points that make strong predictions. Training consistently terminates in one region when using pretraining and in another, non-overlapping region when not using pretraining. Isomap tries to preserve global relative distances (and hence volumes) so the small region corresponding to pretrained models may indicate that the pretraining-based estimator has reduced variance.

reducing test set error. However, unsupervised pretraining can help tasks other than classification, and can act to improve optimization rather than being merely a regularizer. For example, it can improve both train and test reconstruction error for deep autoencoders (Hinton and Salakhutdinov, 2006).

[Erhan et al. \(2010\)](#) performed many experiments to explain several successes of unsupervised pretraining. Both improvements to training error and improvements to test error may be explained in terms of unsupervised pretraining taking the parameters into a region that would otherwise be inaccessible. Neural network training is non-deterministic, and converges to a different function every time it is run. Training may halt at a point where the gradient becomes small, a point where early stopping ends training to prevent overfitting, or at a point where the gradient is large but it is difficult to find a downhill step due to problems such as stochasticity or poor conditioning of the Hessian. Neural networks that receive unsupervised pretraining consistently halt in the same region of function space, while neural networks without pretraining consistently halt in another region. See Fig. 15.1 for a visualization of this phenomenon. The region where pretrained networks arrive is smaller, suggesting that pretraining reduces the variance of the estimation process, which can in turn reduce the risk of severe over-fitting. In other words, unsupervised pretraining initializes neural network parameters into a region that they do not escape, and the results following this initialization are more consistent and less likely to be very bad than without this initialization.

[Erhan et al. \(2010\)](#) also provide some answers as to **when** pretraining works best—the mean and variance of the test error were most reduced by pretraining for deeper networks. Keep in mind that these experiments were performed before the invention and popularization of modern techniques for training very deep networks (rectified linear units, dropout and batch normalization) so less is known about the effect of unsupervised pretraining in conjunction with contemporary approaches.

An important question is how unsupervised pretraining can act as a regularizer. One hypothesis is that pretraining encourages the learning algorithm to discover features that relate to the underlying causes that generate the observed data. This is an important idea motivating many other algorithms besides unsupervised pretraining, and is described further in Sec. 15.3.

Compared to other ways of incorporating this belief by using unsupervised learning, unsupervised pretraining has the disadvantage that it operates with two separate training phases. One reason that these two training phases are disadvantageous is that there is not a single hyperparameter that predictably reduces or increases the strength of the regularization arising from the unsupervised pretraining. Instead, there are very many hyperparameters, whose effect may be

measured after the fact but is often difficult to predict ahead of time. When we perform unsupervised and supervised learning simultaneously, instead of using the pretraining strategy, there is a single hyperparameter, usually a coefficient attached to the unsupervised cost, that determines how strongly the unsupervised objective will regularize the supervised model. One can always predictably obtain less regularization by decreasing this coefficient. In the case of unsupervised pretraining, there is not a way of flexibly adapting the strength of the regularization—either the supervised model is initialized to pretrained parameters, or it is not.

Another disadvantage of having two separate training phases is that each phase has its own hyperparameters. The performance of the second phase usually cannot be predicted during the first phase, so there is a long delay between proposing hyperparameters for the first phase and being able to update them using feedback from the second phase. The most principled approach is to use validation set error in the supervised phase in order to select the hyperparameters of the pretraining phase, as discussed in [Larochelle et al. \(2009\)](#). In practice, some hyperparameters, like the number of pretraining iterations, are more conveniently set during the pretraining phase, using early stopping on the unsupervised objective, which is not ideal but computationally much cheaper than using the supervised objective.

Today, unsupervised pretraining has been largely abandoned, except in the field of natural language processing, where the natural representation of words as one-hot vectors conveys no similarity information and where very large unlabeled sets are available. In that case, the advantage of pretraining is that one can pretrain once on a huge unlabeled set (for example with a corpus containing billions of words), learn a good representation (typically of words, but also of sentences), and then use this representation or fine-tune it for a supervised task for which the training set contains substantially fewer examples. This approach was pioneered by [Collobert and Weston \(2008b\)](#), [Turian et al. \(2010\)](#), and [Collobert et al. \(2011a\)](#) and remains in common use today.

Deep learning techniques based on supervised learning, regularized with dropout or batch normalization, are able to achieve human-level performance on very many tasks, but only with extremely large labeled datasets. These same techniques outperform unsupervised pretraining on medium-sized datasets such as CIFAR-10 and MNIST, which have roughly 5,000 labeled examples per class. On extremely small datasets, such as the alternative splicing dataset, Bayesian methods outperform methods based on unsupervised pretraining ([Srivastava, 2013](#)). For these reasons, the popularity of unsupervised pretraining has declined. Nevertheless, unsupervised pretraining remains an important milestone in the history of deep learning research and continues to influence contemporary approaches. The idea of

pretraining has been generalized to *supervised pretraining* discussed in Sec. 8.7.4, as a very common approach for transfer learning. Supervised pretraining for transfer learning is popular (Oquab *et al.*, 2014; Yosinski *et al.*, 2014) for use with convolutional networks pretrained on the ImageNet dataset. Practitioners publish the parameters of these trained networks for this purpose, just like pretrained word vectors are published for natural language tasks (Collobert *et al.*, 2011a; Mikolov *et al.*, 2013a).

## 15.2 Transfer Learning and Domain Adaptation

Transfer learning and domain adaptation refer to the situation where what has been learned in one setting (i.e., distribution  $P_1$ ) is exploited to improve generalization in another setting (say distribution  $P_2$ ). This generalizes the idea presented in the previous section, where we transferred representations between an unsupervised learning task and a supervised learning task.

In *transfer learning*, the learner must perform two or more different tasks, but we assume that many of the factors that explain the variations in  $P_1$  are relevant to the variations that need to be captured for learning  $P_2$ . This is typically understood in a supervised learning context, where the input is the same but the target may be of a different nature. For example, we may learn about one set of visual categories, such as cats and dogs, in the first setting, then learn about a different set of visual categories, such as ants and wasps, in the second setting. If there is significantly more data in the first setting (sampled from  $P_1$ ), then that may help to learn representations that are useful to quickly generalize from only very few examples drawn from  $P_2$ . Many visual categories **share** low-level notions of edges and visual shapes, the effects of geometric changes, changes in lighting, etc. In general, transfer learning, multi-task learning (Sec. 7.7), and domain adaptation can be achieved via representation learning when there exist features that are useful for the different settings or tasks, corresponding to underlying factors that appear in more than one setting. This is illustrated in Fig. 7.2, with shared lower layers and task-dependent upper layers.

However, sometimes, what is shared among the different tasks is not the semantics of the input but the semantics of the output. For example, a speech recognition system needs to produce valid sentences at the output layer, but the earlier layers near the input may need to recognize very different versions of the same phonemes or sub-phonemic vocalizations depending on which person is speaking. In cases like these, it makes more sense to share the upper layers (near the output) of the neural network, and have a task-specific preprocessing, as

illustrated in Fig. 15.2.

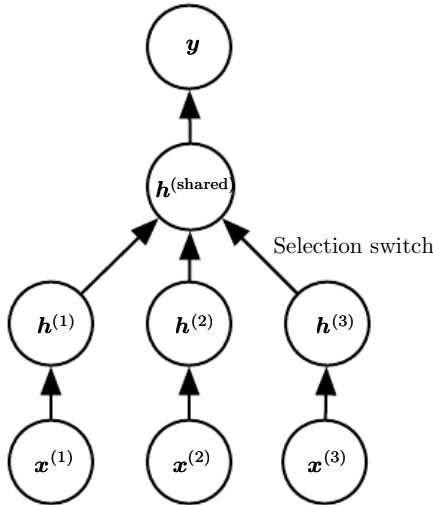


Figure 15.2: Example architecture for multi-task or transfer learning when the output variable  $y$  has the same semantics for all tasks while the input variable  $x$  has a different meaning (and possibly even a different dimension) for each task (or, for example, each user), called  $\mathbf{x}^{(1)}$ ,  $\mathbf{x}^{(2)}$  and  $\mathbf{x}^{(3)}$  for three tasks. The lower levels (up to the selection switch) are task-specific, while the upper levels are shared. The lower levels learn to translate their task-specific input into a generic set of features.

In the related case of *domain adaptation*, the task (and the optimal input-to-output mapping) remains the same between each setting, but the input distribution is slightly different. For example, consider the task of sentiment analysis, which consists of determining whether a comment expresses positive or negative sentiment. Comments posted on the web come from many categories. A domain adaptation scenario can arise when a sentiment predictor trained on customer reviews of media content such as books, videos and music is later used to analyze comments about consumer electronics such as televisions or smartphones. One can imagine that there is an underlying function that tells whether any statement is positive, neutral or negative, but of course the vocabulary and style may vary from one domain to another, making it more difficult to generalize across domains. Simple unsupervised pretraining (with denoising autoencoders) has been found to be very successful for sentiment analysis with domain adaptation (Glorot *et al.*, 2011b).

A related problem is that of *concept drift*, which we can view as a form of transfer learning due to gradual changes in the data distribution over time. Both concept drift and transfer learning can be viewed as particular forms of multi-task learning.

While the phrase “multi-task learning” typically refers to supervised learning tasks, the more general notion of transfer learning is applicable to unsupervised learning and reinforcement learning as well.

In all of these cases, the objective is to take advantage of data from the first setting to extract information that may be useful when learning or even when directly making predictions in the second setting. The core idea of representation learning is that the same representation may be useful in both settings. Using the same representation in both settings allows the representation to benefit from the training data that is available for both tasks.

As mentioned before, unsupervised deep learning for transfer learning has found success in some machine learning competitions (Mesnil *et al.*, 2011; Goodfellow *et al.*, 2011). In the first of these competitions, the experimental setup is the following. Each participant is first given a dataset from the first setting (from distribution  $P_1$ ), illustrating examples of some set of categories. The participants must use this to learn a good feature space (mapping the raw input to some representation), such that when we apply this learned transformation to inputs from the transfer setting (distribution  $P_2$ ), a linear classifier can be trained and generalize well from very few labeled examples. One of the most striking results found in this competition is that as an architecture makes use of deeper and deeper representations (learned in a purely unsupervised way from data collected in the first setting,  $P_1$ ), the learning curve on the new categories of the second (transfer) setting  $P_2$  becomes much better. For deep representations, fewer labeled examples of the transfer tasks are necessary to achieve the apparently asymptotic generalization performance.

Two extreme forms of transfer learning are *one-shot learning* and *zero-shot learning*, sometimes also called *zero-data learning*. Only one labeled example of the transfer task is given for one-shot learning, while no labeled examples are given at all for the zero-shot learning task.

One-shot learning (Fei-Fei *et al.*, 2006) is possible because the representation learns to cleanly separate the underlying classes during the first stage. During the transfer learning stage, only one labeled example is needed to infer the label of many possible test examples that all cluster around the same point in representation space. This works to the extent that the factors of variation corresponding to these invariances have been cleanly separated from the other factors, in the learned representation space, and we have somehow learned which factors do and do not matter when discriminating objects of certain categories.

As an example of a zero-shot learning setting, consider the problem of having a learner read a large collection of text and then solve object recognition problems.

It may be possible to recognize a specific object class even without having seen an image of that object, if the text describes the object well enough. For example, having read that a cat has four legs and pointy ears, the learner might be able to guess that an image is a cat, without having seen a cat before.

Zero-data learning (Larochelle *et al.*, 2008) and zero-shot learning (Palatucci *et al.*, 2009; Socher *et al.*, 2013b) are only possible because additional information has been exploited during training. We can think of the zero-data learning scenario as including three random variables: the traditional inputs  $\mathbf{x}$ , the traditional outputs or targets  $\mathbf{y}$ , and an additional random variable describing the task,  $T$ . The model is trained to estimate the conditional distribution  $p(\mathbf{y} \mid \mathbf{x}, T)$  where  $T$  is a description of the task we wish the model to perform. In our example of recognizing cats after having read about cats, the output is a binary variable  $y$  with  $y = 1$  indicating “yes” and  $y = 0$  indicating “no.” The task variable  $T$  then represents questions to be answered such as “Is there a cat in this image?” If we have a training set containing unsupervised examples of objects that live in the same space as  $T$ , we may be able to infer the meaning of unseen instances of  $T$ . In our example of recognizing cats without having seen an image of the cat, it is important that we have had unlabeled text data containing sentences such as “cats have four legs” or “cats have pointy ears.”

Zero-shot learning requires  $T$  to be represented in a way that allows some sort of generalization. For example,  $T$  cannot be just a one-hot code indicating an object category. Socher *et al.* (2013b) provide instead a distributed representation of object categories by using a learned word embedding for the word associated with each category.

A similar phenomenon happens in machine translation (Klementiev *et al.*, 2012; Mikolov *et al.*, 2013b; Gouws *et al.*, 2014): we have words in one language, and the relationships between words can be learned from unilingual corpora; on the other hand, we have translated sentences which relate words in one language with words in the other. Even though we may not have labeled examples translating word  $A$  in language  $X$  to word  $B$  in language  $Y$ , we can generalize and guess a translation for word  $A$  because we have learned a distributed representation for words in language  $X$ , a distributed representation for words in language  $Y$ , and created a link (possibly two-way) relating the two spaces, via training examples consisting of matched pairs of sentences in both languages. This transfer will be most successful if all three ingredients (the two representations and the relations between them) are learned jointly.

Zero-shot learning is a particular form of transfer learning. The same principle explains how one can perform *multi-modal learning*, capturing a representation in

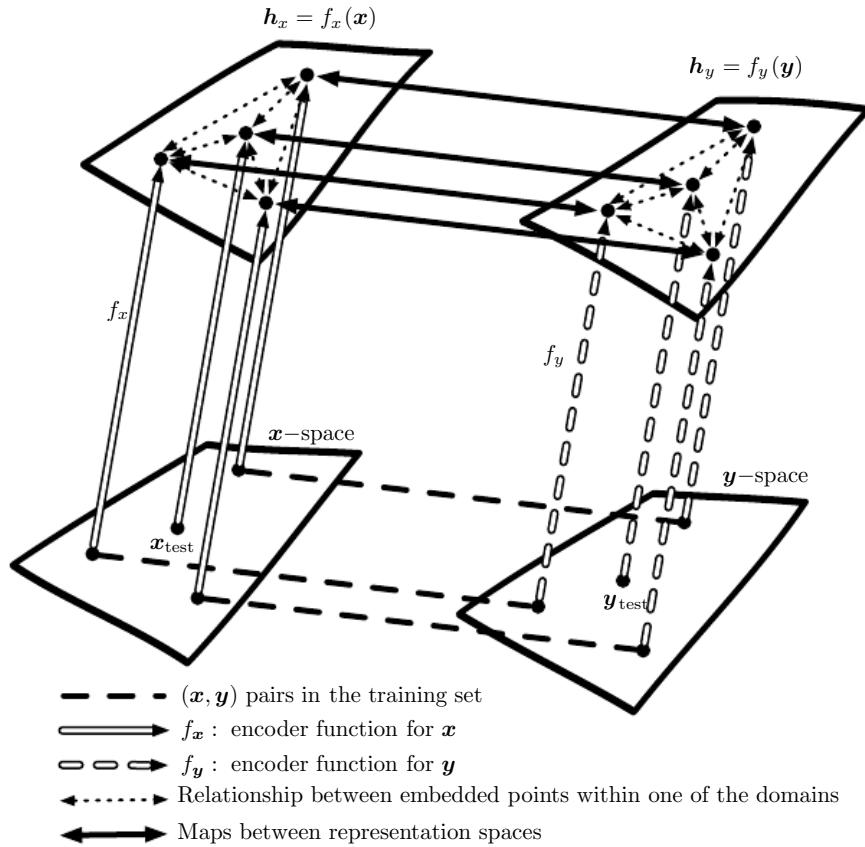


Figure 15.3: Transfer learning between two domains  $\mathbf{x}$  and  $\mathbf{y}$  enables zero-shot learning. Labeled or unlabeled examples of  $\mathbf{x}$  allow one to learn a representation function  $f_x$  and similarly with examples of  $\mathbf{y}$  to learn  $f_y$ . Each application of the  $f_x$  and  $f_y$  functions appears as an upward arrow, with the style of the arrows indicating which function is applied. Distance in  $h_x$  space provides a similarity metric between any pair of points in  $\mathbf{x}$  space that may be more meaningful than distance in  $\mathbf{x}$  space. Likewise, distance in  $h_y$  space provides a similarity metric between any pair of points in  $\mathbf{y}$  space. Both of these similarity functions are indicated with dotted bidirectional arrows. Labeled examples (dashed horizontal lines) are pairs  $(\mathbf{x}, \mathbf{y})$  which allow one to learn a one-way or two-way map (solid bidirectional arrow) between the representations  $f_x(\mathbf{x})$  and the representations  $f_y(\mathbf{y})$  and anchor these representations to each other. Zero-data learning is then enabled as follows. One can associate an image  $\mathbf{x}_{\text{test}}$  to a word  $\mathbf{y}_{\text{test}}$ , even if no image of that word was ever presented, simply because word-representations  $f_y(\mathbf{y}_{\text{test}})$  and image-representations  $f_x(\mathbf{x}_{\text{test}})$  can be related to each other via the maps between representation spaces. It works because, although that image and that word were never paired, their respective feature vectors  $f_x(\mathbf{x}_{\text{test}})$  and  $f_y(\mathbf{y}_{\text{test}})$  have been related to each other. Figure inspired from suggestion by Hrant Khachatrian.

one modality, a representation in the other, and the relationship (in general a joint distribution) between pairs  $(\mathbf{x}, \mathbf{y})$  consisting of one observation  $\mathbf{x}$  in one modality and another observation  $\mathbf{y}$  in the other modality (Srivastava and Salakhutdinov, 2012). By learning all three sets of parameters (from  $\mathbf{x}$  to its representation, from  $\mathbf{y}$  to its representation, and the relationship between the two representations), concepts in one representation are anchored in the other, and vice-versa, allowing one to meaningfully generalize to new pairs. The procedure is illustrated in Fig. 15.3.

### 15.3 Semi-Supervised Disentangling of Causal Factors

An important question about representation learning is “what makes one representation better than another?” One hypothesis is that an ideal representation is one in which the features within the representation correspond to the underlying causes of the observed data, with separate features or directions in feature space corresponding to different causes, so that the representation disentangles the causes from one another. This hypothesis motivates approaches in which we first seek a good representation for  $p(\mathbf{x})$ . Such a representation may also be a good representation for computing  $p(\mathbf{y} | \mathbf{x})$  if  $\mathbf{y}$  is among the most salient causes of  $\mathbf{x}$ . This idea has guided a large amount of deep learning research since at least the 1990s (Becker and Hinton, 1992; Hinton and Sejnowski, 1999), in more detail. For other arguments about when semi-supervised learning can outperform pure supervised learning, we refer the reader to Sec. 1.2 of Chapelle *et al.* (2006).

In other approaches to representation learning, we have often been concerned with a representation that is easy to model—for example, one whose entries are sparse, or independent from each other. A representation that cleanly separates the underlying causal factors may not necessarily be one that is easy to model. However, a further part of the hypothesis motivating semi-supervised learning via unsupervised representation learning is that for many AI tasks, these two properties coincide: once we are able to obtain the underlying explanations for what we observe, it generally becomes easy to isolate individual attributes from the others. Specifically, if a representation  $\mathbf{h}$  represents many of the underlying causes of the observed  $\mathbf{x}$ , and the outputs  $\mathbf{y}$  are among the most salient causes, then it is easy to predict  $\mathbf{y}$  from  $\mathbf{h}$ .

First, let us see how semi-supervised learning can fail because unsupervised learning of  $p(\mathbf{x})$  is of no help to learn  $p(\mathbf{y} | \mathbf{x})$ . Consider for example the case where  $p(\mathbf{x})$  is uniformly distributed and we want to learn  $f(\mathbf{x}) = \mathbb{E}[\mathbf{y} | \mathbf{x}]$ . Clearly, observing a training set of  $\mathbf{x}$  values alone gives us no information about  $p(\mathbf{y} | \mathbf{x})$ .

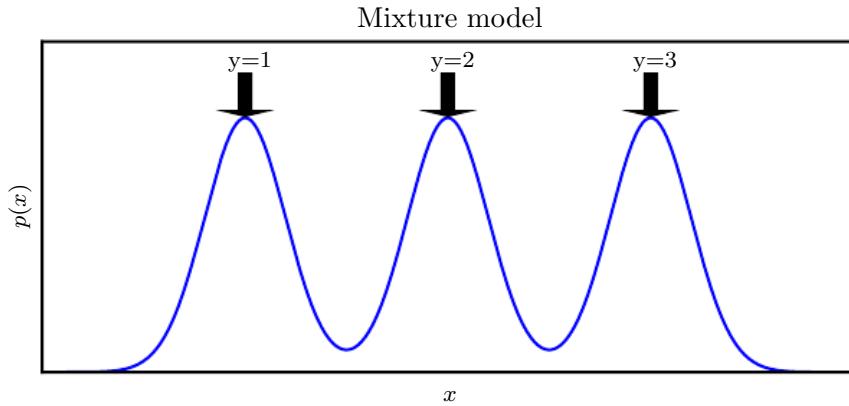


Figure 15.4: Example of a density over  $x$  that is a mixture over three components. The component identity is an underlying explanatory factor,  $y$ . Because the mixture components (e.g., natural object classes in image data) are statistically salient, just modeling  $p(x)$  in an unsupervised way with no labeled example already reveals the factor  $y$ .

Next, let us see a simple example of how semi-supervised learning can succeed. Consider the situation where  $\mathbf{x}$  arises from a mixture, with one mixture component per value of  $\mathbf{y}$ , as illustrated in Fig. 15.4. If the mixture components are well-separated, then modeling  $p(\mathbf{x})$  reveals precisely where each component is, and a single labeled example of each class will then be enough to perfectly learn  $p(\mathbf{y} | \mathbf{x})$ . But more generally, what could make  $p(\mathbf{y} | \mathbf{x})$  and  $p(\mathbf{x})$  be tied together?

If  $\mathbf{y}$  is closely associated with one of the causal factors of  $\mathbf{x}$ , then  $p(\mathbf{x})$  and  $p(\mathbf{y} | \mathbf{x})$  will be strongly tied, and unsupervised representation learning that tries to disentangle the underlying factors of variation is likely to be useful as a semi-supervised learning strategy.

Consider the assumption that  $\mathbf{y}$  is one of the causal factors of  $\mathbf{x}$ , and let  $\mathbf{h}$  represent all those factors. The true generative process can be conceived as structured according to this directed graphical model, with  $\mathbf{h}$  as the parent of  $\mathbf{x}$ :

$$p(\mathbf{h}, \mathbf{x}) = p(\mathbf{x} | \mathbf{h})p(\mathbf{h}). \quad (15.1)$$

As a consequence, the data has marginal probability

$$p(\mathbf{x}) = \mathbb{E}_{\mathbf{h}} p(\mathbf{x} | \mathbf{h}). \quad (15.2)$$

From this straightforward observation, we conclude that the best possible model of  $\mathbf{x}$  (from a generalization point of view) is the one that uncovers the above “true”

structure, with  $\mathbf{h}$  as a latent variable that explains the observed variations in  $\mathbf{x}$ . The “ideal” representation learning discussed above should thus recover these latent factors. If  $\mathbf{y}$  is one of these (or closely related to one of them), then it will be very easy to learn to predict  $\mathbf{y}$  from such a representation. We also see that the conditional distribution of  $\mathbf{y}$  given  $\mathbf{x}$  is tied by Bayes rule to the components in the above equation:

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid \mathbf{y})p(\mathbf{y})}{p(\mathbf{x})}. \quad (15.3)$$

Thus the marginal  $p(\mathbf{x})$  is intimately tied to the conditional  $p(\mathbf{y} \mid \mathbf{x})$  and knowledge of the structure of the former should be helpful to learn the latter. Therefore, in situations respecting these assumptions, semi-supervised learning should improve performance.

An important research problem regards the fact that most observations are formed by an extremely large number of underlying causes. Suppose  $\mathbf{y} = \mathbf{h}_i$ , but the unsupervised learner does not know which  $\mathbf{h}_i$ . The brute force solution is for an unsupervised learner to learn a representation that captures **all** the reasonably salient generative factors  $\mathbf{h}_j$  and disentangles them from each other, thus making it easy to predict  $\mathbf{y}$  from  $\mathbf{h}$ , regardless of which  $\mathbf{h}_i$  is associated with  $\mathbf{y}$ .

In practice, the brute force solution is not feasible because it is not possible to capture all or most of the factors of variation that influence an observation. For example, in a visual scene, should the representation always encode all of the smallest objects in the background? It is a well-documented psychological phenomenon that human beings fail to perceive changes in their environment that are not immediately relevant to the task they are performing—see, e.g., [Simons and Levin \(1998\)](#). An important research frontier in semi-supervised learning is determining **what** to encode in each situation. Currently, two of the main strategies for dealing with a large number of underlying causes are to use a supervised learning signal at the same time as the unsupervised learning signal so that the model will choose to capture the most relevant factors of variation, or to use much larger representations if using purely unsupervised learning.

An emerging strategy for unsupervised learning is to modify the definition of which underlying causes are most salient. Historically, autoencoders and generative models have been trained to optimize a fixed criterion, often similar to mean squared error. These fixed criteria determine which causes are considered salient. For example, mean squared error applied to the pixels of an image implicitly specifies that an underlying cause is only salient if it significantly changes the brightness of a large number of pixels. This can be problematic if the task we wish to solve involves interacting with small objects. See Fig. 15.5 for an example

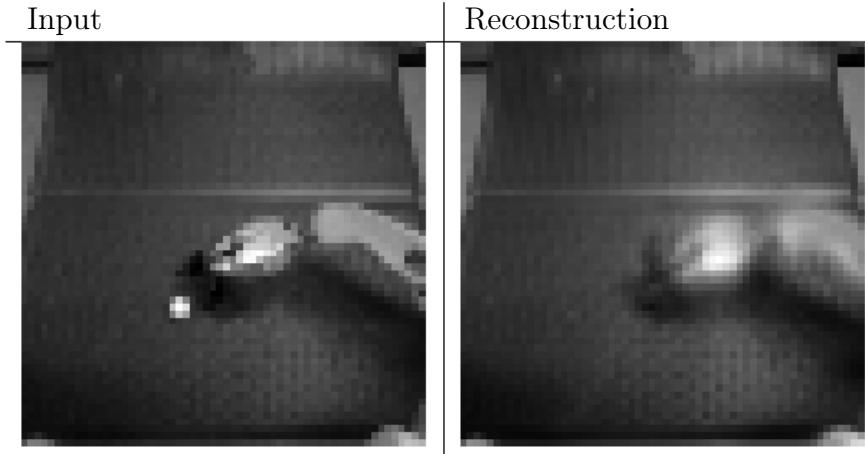


Figure 15.5: An autoencoder trained with mean squared error for a robotics task has failed to reconstruct a ping pong ball. The existence of the ping pong ball and all of its spatial coordinates are important underlying causal factors that generate the image and are relevant to the robotics task. Unfortunately, the autoencoder has limited capacity, and the training with mean squared error did not identify the ping pong ball as being salient enough to encode. Images graciously provided by Chelsea Finn.

of a robotics task in which an autoencoder has failed to learn to encode a small ping pong ball. This same robot is capable of successfully interacting with larger objects, such as baseballs, which are more salient according to mean squared error.

Other definitions of salience are possible. For example, if a group of pixels follow a highly recognizable pattern, even if that pattern does not involve extreme brightness or darkness, then that pattern could be considered extremely salient. One way to implement such a definition of salience is to use a recently developed approach called *generative adversarial networks* (Goodfellow *et al.*, 2014c). In this approach, a generative model is trained to fool a feedforward classifier. The feedforward classifier attempts to recognize all samples from the generative model as being fake, and all samples from the training set as being real. In this framework, any structured pattern that the feedforward network can recognize is highly salient. The generative adversarial network will be described in more detail in Sec. 20.10.4. For the purposes of the present discussion, it is sufficient to understand that they **learn** how to determine what is salient. Lotter *et al.* (2015) showed that models trained to generate images of human heads will often neglect to generate the ears when trained with mean squared error, but will successfully generate the ears when trained with the adversarial framework. Because the ears are not extremely bright or dark compared to the surrounding skin, they are not especially salient according to mean squared error loss, but their highly recognizable shape and consistent

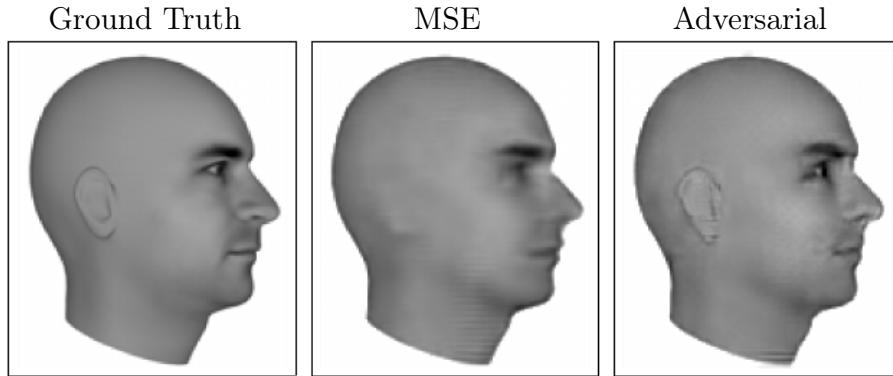


Figure 15.6: Predictive generative networks provide an example of the importance of learning which features are salient. In this example, the predictive generative network has been trained to predict the appearance of a 3-D model of a human head at a specific viewing angle. (*Left*) Ground truth. This is the correct image, that the network should emit. (*Center*) Image produced by a predictive generative network trained with mean squared error alone. Because the ears do not cause an extreme difference in brightness compared to the neighboring skin, they were not sufficiently salient for the model to learn to represent them. (*Right*) Image produced by a model trained with a combination of mean squared error and adversarial loss. Using this learned cost function, the ears are salient because they follow a predictable pattern. Learning which underlying causes are important and relevant enough to model is an important active area of research. Figures graciously provided by [Lotter et al. \(2015\)](#).

position means that a feedforward network can easily learn to detect them, making them highly salient under the generative adversarial framework. See Fig. 15.6 for example images. Generative adversarial networks are only one step toward determining which factors should be represented. We expect that future research will discover better ways of determining which factors to represent, and develop mechanisms for representing different factors depending on the task.

A benefit of learning the underlying causal factors, as pointed out by [Schölkopf et al. \(2012\)](#), is that if the true generative process has  $\mathbf{x}$  as an effect and  $\mathbf{y}$  as a cause, then modeling  $p(\mathbf{x} \mid \mathbf{y})$  is robust to changes in  $p(\mathbf{y})$ . If the cause-effect relationship was reversed, this would not be true, since by Bayes rule,  $p(\mathbf{x} \mid \mathbf{y})$  would be sensitive to changes in  $p(\mathbf{y})$ . Very often, when we consider changes in distribution due to different domains, temporal non-stationarity, or changes in the nature of the task, **the causal mechanisms remain invariant** (“the laws of the universe are constant”) while the marginal distribution over the underlying causes can change. Hence, better generalization and robustness to all kinds of changes can be expected via learning a generative model that attempts to recover

the causal factors  $\mathbf{h}$  and  $p(\mathbf{x} \mid \mathbf{h})$ .

## 15.4 Distributed Representation

Distributed representations of concepts—representations composed of many elements that can be set separately from each other—are one of the most important tools for representation learning. Distributed representations are powerful because they can use  $n$  features with  $k$  values to describe  $k^n$  different concepts. As we have seen throughout this book, both neural networks with multiple hidden units and probabilistic models with multiple latent variables make use of the strategy of distributed representation. We now introduce an additional observation. Many deep learning algorithms are motivated by the assumption that the hidden units can learn to represent the underlying causal factors that explain the data, as discussed in Sec. 15.3. Distributed representations are natural for this approach, because each direction in representation space can correspond to the value of a different underlying configuration variable.

An example of a distributed representation is a vector of  $n$  binary features, which can take  $2^n$  configurations, each potentially corresponding to a different region in input space, as illustrated in Fig. 15.7. This can be compared with a *symbolic representation*, where the input is associated with a single symbol or category. If there are  $n$  symbols in the dictionary, one can imagine  $n$  feature detectors, each corresponding to the detection of the presence of the associated category. In that case only  $n$  different configurations of the representation space are possible, carving  $n$  different regions in input space, as illustrated in Fig. 15.8. Such a symbolic representation is also called a one-hot representation, since it can be captured by a binary vector with  $n$  bits that are mutually exclusive (only one of them can be active). A symbolic representation is a specific example of the broader class of non-distributed representations, which are representations that may contain many entries but without significant meaningful separate control over each entry.

Examples of learning algorithms based on non-distributed representations include:

- Clustering methods, including the  $k$ -means algorithm: each input point is assigned to exactly one cluster.
- $k$ -nearest neighbors algorithms: one or a few templates or prototype examples are associated with a given input. In the case of  $k > 1$ , there are multiple

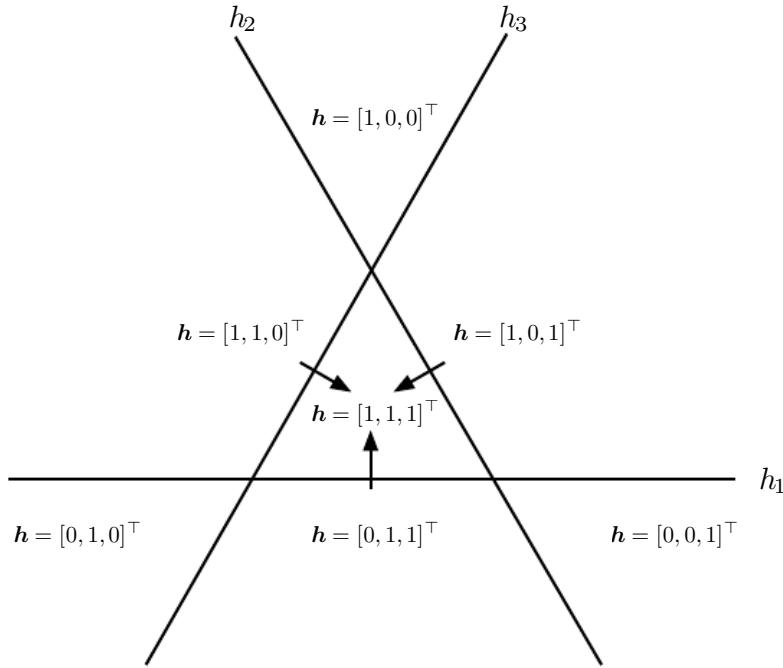


Figure 15.7: Illustration of how a learning algorithm based on a distributed representation breaks up the input space into regions. In this example, there are three binary features  $h_1$ ,  $h_2$ , and  $h_3$ . Each feature is defined by thresholding the output of a learned, linear transformation. Each feature divides  $\mathbb{R}^2$  into two half-planes. Let  $h_i^+$  be the set of input points for which  $h_i = 1$  and  $h_i^-$  be the set of input points for which  $h_i = 0$ . In this illustration, each line represents the decision boundary for one  $h_i$ , with the corresponding arrow pointing to the  $h_i^+$  side of the boundary. The representation as a whole takes on a unique value at each possible intersection of these half-planes. For example, the representation value  $[1, 1, 1]^\top$  corresponds to the region  $h_1^+ \cap h_2^+ \cap h_3^+$ . Compare this to the non-distributed representations in Fig. 15.8. In the general case of  $d$  input dimensions, a distributed representation divides  $\mathbb{R}^d$  by intersecting half-spaces rather than half-planes. The distributed representation with  $n$  features assigns unique codes to  $O(n^d)$  different regions, while the nearest neighbor algorithm with  $n$  examples assigns unique codes to only  $n$  regions. The distributed representation is thus able to distinguish exponentially many more regions than the non-distributed one. Keep in mind that not all  $\mathbf{h}$  values are feasible (there is no  $\mathbf{h} = \mathbf{0}$  in this example) and that a linear classifier on top of the distributed representation is not able to assign different class identities to every neighboring region; even a deep linear-threshold network has a VC dimension of only  $O(w \log w)$  where  $w$  is the number of weights (Sontag, 1998). The combination of a powerful representation layer and a weak classifier layer can be a strong regularizer; a classifier trying to learn the concept of “person” versus “not a person” does not need to assign a different class to an input represented as “woman with glasses” than it assigns to an input represented as “man without glasses.” This capacity constraint encourages each classifier to focus on few  $h_i$  and encourages  $\mathbf{h}$  to learn to represent the classes in a linearly separable way.

values describing each input, but they can not be controlled separately from each other, so this does not qualify as a true distributed representation.

- Decision trees: only one leaf (and the nodes on the path from root to leaf) is activated when an input is given.
- Gaussian mixtures and mixtures of experts: the templates (cluster centers) or experts are now associated with a **degree** of activation. As with the  $k$ -nearest neighbors algorithm, each input is represented with multiple values, but those values cannot readily be controlled separately from each other.
- Kernel machines with a Gaussian kernel (or other similarly local kernel): although the degree of activation of each “support vector” or template example is now continuous-valued, the same issue arises as with Gaussian mixtures.
- Language or translation models based on  $n$ -grams. The set of contexts (sequences of symbols) is partitioned according to a tree structure of suffixes. A leaf may correspond to the last two words being  $w_1$  and  $w_2$ , for example. Separate parameters are estimated for each leaf of the tree (with some sharing being possible).

For some of these non-distributed algorithms, the output is not constant by parts but instead interpolates between neighboring regions. The relationship between the number of parameters (or examples) and the number of regions they can define remains linear.

An important related concept that distinguishes a distributed representation from a symbolic one is that **generalization arises due to shared attributes** between different concepts. As pure symbols, “cat” and “dog” are as far from each other as any other two symbols. However, if one associates them with a meaningful distributed representation, then many of the things that can be said about cats can generalize to dogs and vice-versa. For example, our distributed representation may contain entries such as “`has_fur`” or “`number_of_legs`” that have the same value for the embedding of both “`cat`” and “`dog`. Neural language models that operate on distributed representations of words generalize much better than other models that operate directly on one-hot representations of words, as discussed in Sec. 12.4. Distributed representations induce a rich *similarity space*, in which semantically close concepts (or inputs) are close in distance, a property that is absent from purely symbolic representations.

When and why can there be a statistical advantage from using a distributed representation as part of a learning algorithm? Distributed representations can

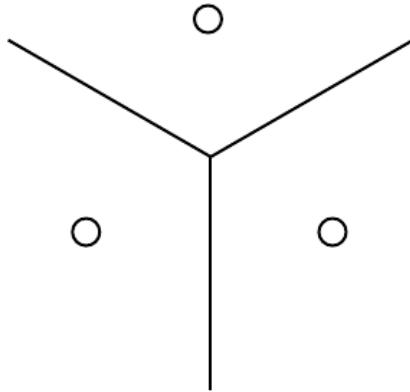


Figure 15.8: Illustration of how the nearest neighbor algorithm breaks up the input space into different regions. The nearest neighbor algorithm provides an example of a learning algorithm based on a non-distributed representation. Different non-distributed algorithms may have different geometry, but they typically break the input space into regions, **with a separate set of parameters for each region**. The advantage of a non-distributed approach is that, given enough parameters, it can fit the training set without solving a difficult optimization algorithm, because it is straightforward to choose a different output **independently** for each region. The disadvantage is that such non-distributed models generalize only locally via the smoothness prior, making it difficult to learn a complicated function with more peaks and troughs than the available number of examples. Contrast this with a distributed representation, Fig. 15.7.

have a statistical advantage when an apparently complicated structure can be compactly represented using a small number of parameters. Some traditional non-distributed learning algorithms generalize only due to the smoothness assumption, which states that if  $u \approx v$ , then the target function  $f$  to be learned has the property that  $f(u) \approx f(v)$ , in general. There are many ways of formalizing such an assumption, but the end result is that if we have an example  $(x, y)$  for which we know that  $f(x) \approx y$ , then we choose an estimator  $\hat{f}$  that approximately satisfies these constraints while changing as little as possible when we move to a nearby input  $x + \epsilon$ . This assumption is clearly very useful, but it suffers from the curse of dimensionality: in order to learn a target function that increases and decreases many times in many different regions,<sup>1</sup> we may need a number of examples that is at least as large as the number of distinguishable regions. One can think of each of these regions as a category or symbol: by having a separate degree of freedom for each symbol (or region), we can learn an arbitrary decoder mapping from symbol to value. However, this does not allow us to generalize to new symbols for new regions.

If we are lucky, there may be some regularity in the target function, besides being smooth. For example, a convolutional network with max-pooling can recognize an object regardless of its location in the image, even though spatial translation of the object may not correspond to smooth transformations in the input space.

Let us examine a special case of a distributed representation learning algorithm, that extracts binary features by thresholding linear functions of the input. Each binary feature in this representation divides  $\mathbb{R}^d$  into a pair of half-spaces, as illustrated in Fig. 15.7. The exponentially large number of intersections of  $n$  of the corresponding half-spaces determines how many regions this distributed representation learner can distinguish. How many regions are generated by an arrangement of  $n$  hyperplanes in  $\mathbb{R}^d$ ? By applying a general result concerning the intersection of hyperplanes (Zaslavsky, 1975), one can show (Pascanu *et al.*, 2014b) that the number of regions this binary feature representation can distinguish is

$$\sum_{j=0}^d \binom{n}{j} = O(n^d). \quad (15.4)$$

Therefore, we see a growth that is exponential in the input size and polynomial in the number of hidden units.

---

<sup>1</sup>Potentially, we may want to learn a function whose behavior is distinct in exponentially many regions: in a  $d$ -dimensional space with at least 2 different values to distinguish per dimension, we might want  $f$  to differ in  $2^d$  different regions, requiring  $O(2^d)$  training examples.

This provides a geometric argument to explain the generalization power of distributed representation: with  $O(nd)$  parameters (for  $n$  linear-threshold features in  $\mathbb{R}^d$ ) we can distinctly represent  $O(n^d)$  regions in input space. If instead we made no assumption at all about the data, and used a representation with one unique symbol for each region, and separate parameters for each symbol to recognize its corresponding portion of  $\mathbb{R}^d$ , then specifying  $O(n^d)$  regions would require  $O(n^d)$  examples. More generally, the argument in favor of the distributed representation could be extended to the case where instead of using linear threshold units we use nonlinear, possibly continuous, feature extractors for each of the attributes in the distributed representation. The argument in this case is that if a parametric transformation with  $k$  parameters can learn about  $r$  regions in input space, with  $k \ll r$ , and if obtaining such a representation was useful to the task of interest, then we could potentially generalize much better in this way than in a non-distributed setting where we would need  $O(r)$  examples to obtain the same features and associated partitioning of the input space into  $r$  regions. Using fewer parameters to represent the model means that we have fewer parameters to fit, and thus require far fewer training examples to generalize well.

A further part of the argument for why models based on distributed representations generalize well is that their capacity remains limited despite being able to distinctly encode so many different regions. For example, the VC dimension of a neural network of linear threshold units is only  $O(w \log w)$ , where  $w$  is the number of weights (Sontag, 1998). This limitation arises because, while we can assign very many unique codes to representation space, we cannot use absolutely all of the code space, nor can we learn arbitrary functions mapping from the representation space  $\mathbf{h}$  to the output  $\mathbf{y}$  using a linear classifier. The use of a distributed representation combined with a linear classifier thus expresses a prior belief that the classes to be recognized are linearly separable as a function of the underlying causal factors captured by  $\mathbf{h}$ . We will typically want to learn categories such as the set of all images of all green objects or the set of all images of cars, but not categories that require nonlinear, XOR logic. For example, we typically do not want to partition the data into the set of all red cars and green trucks as one class and the set of all green cars and red trucks as another class.

The ideas discussed so far have been abstract, but they may be experimentally validated. Zhou *et al.* (2015) find that hidden units in a deep convolutional network trained on the ImageNet and Places benchmark datasets learn features that are very often interpretable, corresponding to a label that humans would naturally assign. In practice it is certainly not always the case that hidden units learn something that has a simple linguistic name, but it is interesting to see this emerge near the top levels of the best computer vision deep networks. What such

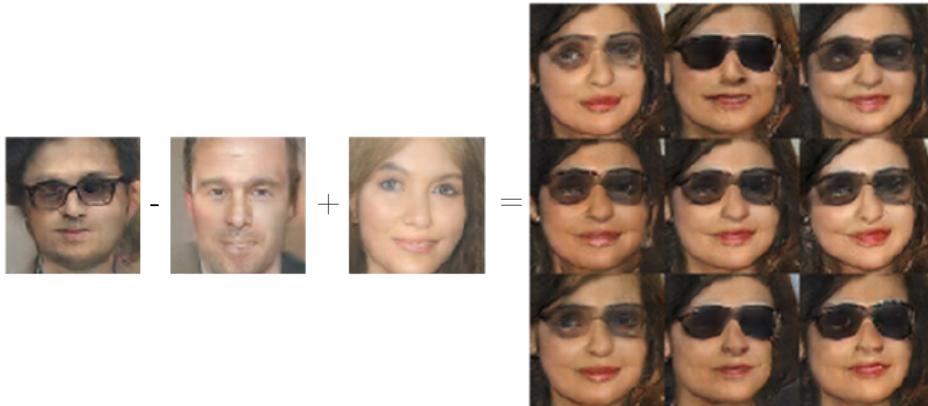


Figure 15.9: A generative model has learned a distributed representation that disentangles the concept of gender from the concept of wearing glasses. If we begin with the representation of the concept of a man with glasses, then subtract the vector representing the concept of a man without glasses, and finally add the vector representing the concept of a woman without glasses, we obtain the vector representing the concept of a woman with glasses. The generative model correctly decodes all of these representation vectors to images that may be recognized as belonging to the correct class. Images reproduced with permission from [Radford et al. \(2015\)](#).

features have in common is that one could imagine **learning about each of them without having to see all the configurations of all the others**. [Radford et al. \(2015\)](#) demonstrated that a generative model can learn a representation of images of faces, with separate directions in representation space capturing different underlying factors of variation. Fig. 15.9 demonstrates that one direction in representation space corresponds to whether the person is male or female, while another corresponds to whether the person is wearing glasses. These features were discovered automatically, not fixed a priori. There is no need to have labels for the hidden unit classifiers: gradient descent on an objective function of interest naturally learns semantically interesting features, so long as the task requires such features. We can learn about the distinction between male and female, or about the presence or absence of glasses, without having to characterize all of the configurations of the  $n - 1$  other features by examples covering all of these combinations of values. This form of statistical separability is what allows one to generalize to new configurations of a person's features that have never been seen during training.

## 15.5 Exponential Gains from Depth

We have seen in Sec. 6.4.1 that multilayer perceptrons are universal approximators, and that some functions can be represented by exponentially smaller deep networks compared to shallow networks. This decrease in model size leads to improved statistical efficiency. In this section, we describe how similar results apply more generally to other kinds of models with distributed hidden representations.

In Sec. 15.4, we saw an example of a generative model that learned about the explanatory factors underlying images of faces, including the person’s gender and whether they are wearing glasses. The generative model that accomplished this task was based on a deep neural network. It would not be reasonable to expect a shallow network, such as a linear network, to learn the complicated relationship between these abstract explanatory factors and the pixels in the image. In this and other AI tasks, the factors that can be chosen almost independently in order to generate data are more likely to be very high-level and related in highly nonlinear ways to the input. We argue that this demands **deep** distributed representations, where the higher level features (seen as functions of the input) or factors (seen as generative causes) are obtained through the composition of many nonlinearities.

It has been proven in many different settings that organizing computation through the composition of many nonlinearities and a hierarchy of reused features can give an exponential boost to statistical efficiency, on top of the exponential boost given by using a distributed representation. Many kinds of networks (e.g., with saturating nonlinearities, Boolean gates, sum/products, or RBF units) with a single hidden layer can be shown to be universal approximators. A model family that is a universal approximator can approximate a large class of functions (including all continuous functions) up to any non-zero tolerance level, given enough hidden units. However, the required number of hidden units may be very large. Theoretical results concerning the expressive power of deep architectures state that there are families of functions that can be represented efficiently by an architecture of depth  $k$ , but would require an exponential number of hidden units (with respect to the input size) with insufficient depth (depth 2 or depth  $k - 1$ ).

In Sec. 6.4.1, we saw that deterministic feedforward networks are universal approximators of functions. Many structured probabilistic models with a single hidden layer of latent variables, including restricted Boltzmann machines and deep belief networks, are universal approximators of probability distributions (Le Roux and Bengio, 2008, 2010; Montúfar and Ay, 2011; Montúfar, 2014; Krause *et al.*, 2013).

In Sec. 6.4.1, we saw that a sufficiently deep feedforward network can have an

exponential advantage over a network that is too shallow. Such results can also be obtained for other models such as probabilistic models. One such probabilistic model is the *sum-product network* or *SPN* (Poon and Domingos, 2011). These models use polynomial circuits to compute the probability distribution over a set of random variables. Delalleau and Bengio (2011) showed that there exist probability distributions for which a minimum depth of SPN is required to avoid needing an exponentially large model. Later, Martens and Medabalimi (2014) showed that there are significant differences between every two finite depths of SPN, and that some of the constraints used to make SPNs tractable may limit their representational power.

Another interesting development is a set of theoretical results for the expressive power of families of deep circuits related to convolutional nets, highlighting an exponential advantage for the deep circuit even when the shallow circuit is allowed to only approximate the function computed by the deep circuit (Cohen *et al.*, 2015). By comparison, previous theoretical work made claims regarding only the case where the shallow circuit must exactly replicate particular functions.

## 15.6 Providing Clues to Discover Underlying Causes

To close this chapter, we come back to one of our original questions: what makes one representation better than another? One answer, first introduced in Sec. 15.3, is that an ideal representation is one that disentangles the underlying causal factors of variation that generated the data, especially those factors that are relevant to our applications. Most strategies for representation learning are based on introducing clues that help the learning to find these underlying factors of variations. The clues can help the learner separate these observed factors from the others. Supervised learning provides a very strong clue: a label  $\mathbf{y}$ , presented with each  $\mathbf{x}$ , that usually specifies the value of at least one of the factors of variation directly. More generally, to make use of abundant unlabeled data, representation learning makes use of other, less direct, hints about the underlying factors. These hints take the form of implicit prior beliefs that we, the designers of the learning algorithm, impose in order to guide the learner. Results such as the no free lunch theorem show that regularization strategies are necessary to obtain good generalization. While it is impossible to find a universally superior regularization strategy, one goal of deep learning is to find a set of fairly generic regularization strategies that are applicable to a wide variety of AI tasks, similar to the tasks that people and animals are able to solve.

We provide here a list of these generic regularization strategies. The list is

clearly not exhaustive, but gives some concrete examples of ways that learning algorithms can be encouraged to discover features that correspond to underlying factors. This list was introduced in Sec. 3.1 of Bengio *et al.* (2013d) and has been partially expanded here.

- **Smoothness:** This is the assumption that  $f(\mathbf{x} + \epsilon\mathbf{d}) \approx f(\mathbf{x})$  for unit  $\mathbf{d}$  and small  $\epsilon$ . This assumption allows the learner to generalize from training examples to nearby points in input space. Many machine learning algorithms leverage this idea, but it is insufficient to overcome the curse of dimensionality.
- **Linearity:** Many learning algorithms assume that relationships between some variables are linear. This allows the algorithm to make predictions even very far from the observed data, but can sometimes lead to overly extreme predictions. Most simple machine learning algorithms that do not make the smoothness assumption instead make the linearity assumption. These are in fact different assumptions—linear functions with large weights applied to high-dimensional spaces may not be very smooth. See Goodfellow *et al.* (2014b) for a further discussion of the limitations of the linearity assumption.
- **Multiple explanatory factors:** Many representation learning algorithms are motivated by the assumption that the data is generated by multiple underlying explanatory factors, and that most tasks can be solved easily given the state of each of these factors. Sec. 15.3 describes how this view motivates semi-supervised learning via representation learning. Learning the structure of  $p(\mathbf{x})$  requires learning some of the same features that are useful for modeling  $p(\mathbf{y} | \mathbf{x})$  because both refer to the same underlying explanatory factors. Sec. 15.4 describes how this view motivates the use of distributed representations, with separate directions in representation space corresponding to separate factors of variation.
- **Causal factors:** the model is constructed in such a way that it treats the factors of variation described by the learned representation  $\mathbf{h}$  as the causes of the observed data  $\mathbf{x}$ , and not vice-versa. As discussed in Sec. 15.3, this is advantageous for semi-supervised learning and makes the learned model more robust when the distribution over the underlying causes changes or when we use the model for a new task.
- **Depth, or a hierarchical organization of explanatory factors:** High-level, abstract concepts can be defined in terms of simple concepts, forming a hierarchy. From another point of view, the use of a deep architecture expresses our belief that the task should be accomplished via a multi-step program,

with each step referring back to the output of the processing accomplished via previous steps.

- **Shared factors across tasks:** In the context where we have many tasks, corresponding to different  $y_i$  variables sharing the same input  $\mathbf{x}$  or where each task is associated with a subset or a function  $f^{(i)}(\mathbf{x})$  of a global input  $\mathbf{x}$ , the assumption is that each  $y_i$  is associated with a different subset from a common pool of relevant factors  $\mathbf{h}$ . Because these subsets overlap, learning all the  $P(y_i | \mathbf{x})$  via a shared intermediate representation  $P(\mathbf{h} | \mathbf{x})$  allows sharing of statistical strength between the tasks.
- **Manifolds:** Probability mass concentrates, and the regions in which it concentrates are locally connected and occupy a tiny volume. In the continuous case, these regions can be approximated by low-dimensional manifolds with a much smaller dimensionality than the original space where the data lives. Many machine learning algorithms behave sensibly only on this manifold ([Goodfellow et al., 2014b](#)). Some machine learning algorithms, especially autoencoders, attempt to explicitly learn the structure of the manifold.
- **Natural clustering:** Many machine learning algorithms assume that each connected manifold in the input space may be assigned to a single class. The data may lie on many disconnected manifolds, but the class remains constant within each one of these. This assumption motivates a variety of learning algorithms, including tangent propagation, double backprop, the manifold tangent classifier and adversarial training.
- **Temporal and spatial coherence:** Slow feature analysis and related algorithms make the assumption that the most important explanatory factors change slowly over time, or at least that it is easier to predict the true underlying explanatory factors than to predict raw observations such as pixel values. See Sec. 13.3 for further description of this approach.
- **Sparsity:** Most features should presumably not be relevant to describing most inputs—there is no need to use a feature that detects elephant trunks when representing an image of a cat. It is therefore reasonable to impose a prior that any feature that can be interpreted as “present” or “absent” should be absent most of the time.
- **Simplicity of Factor Dependencies:** In good high-level representations, the factors are related to each other through simple dependencies. The simplest possible is marginal independence,  $P(\mathbf{h}) = \prod_i P(\mathbf{h}_i)$ , but linear

dependencies or those captured by a shallow autoencoder are also reasonable assumptions. This can be seen in many laws of physics, and is assumed when plugging a linear predictor or a factorized prior on top of a learned representation.

The concept of representation learning ties together all of the many forms of deep learning. Feedforward and recurrent networks, autoencoders and deep probabilistic models all learn and exploit representations. Learning the best possible representation remains an exciting avenue of research.

## Chapter 16

# Structured Probabilistic Models for Deep Learning

Deep learning draws upon many modeling formalisms that researchers can use to guide their design efforts and describe their algorithms. One of these formalisms is the idea of *structured probabilistic models*. We have already discussed structured probabilistic models briefly in Sec. 3.14. That brief presentation was sufficient to understand how to use structured probabilistic models as a language to describe some of the algorithms in Part II. Now, in Part III, structured probabilistic models are a key ingredient of many of the most important research topics in deep learning. In order to prepare to discuss these research ideas, this chapter describes structured probabilistic models in much greater detail. This chapter is intended to be self-contained; the reader does not need to review the earlier introduction before continuing with this chapter.

A structured probabilistic model is a way of describing a probability distribution, using a graph to describe which random variables in the probability distribution interact with each other directly. Here we use “graph” in the graph theory sense—a set of vertices connected to one another by a set of edges. Because the structure of the model is defined by a graph, these models are often also referred to as *graphical models*.

The graphical models research community is large and has developed many different models, training algorithms, and inference algorithms. In this chapter, we provide basic background on some of the most central ideas of graphical models, with an emphasis on the concepts that have proven most useful to the deep learning research community. If you already have a strong background in graphical models, you may wish to skip most of this chapter. However, even a graphical model expert

may benefit from reading the final section of this chapter, Sec. 16.7, in which we highlight some of the unique ways that graphical models are used for deep learning algorithms. Deep learning practitioners tend to use very different model structures, learning algorithms and inference procedures than are commonly used by the rest of the graphical models research community. In this chapter, we identify these differences in preferences and explain the reasons for them.

In this chapter we first describe the challenges of building large-scale probabilistic models. Next, we describe how to use a graph to describe the structure of a probability distribution. While this approach allows us to overcome many challenges, it is not without its own complications. One of the major difficulties in graphical modeling is understanding which variables need to be able to interact directly, i.e., which graph structures are most suitable for a given problem. We outline two approaches to resolving this difficulty by learning about the dependencies in Sec. 16.5. Finally, we close with a discussion of the unique emphasis that deep learning practitioners place on specific approaches to graphical modeling in Sec. 16.7.

## 16.1 The Challenge of Unstructured Modeling

The goal of deep learning is to scale machine learning to the kinds of challenges needed to solve artificial intelligence. This means being able to understand high-dimensional data with rich structure. For example, we would like AI algorithms to be able to understand natural images,<sup>1</sup> audio waveforms representing speech, and documents containing multiple words and punctuation characters.

Classification algorithms can take an input from such a rich high-dimensional distribution and summarize it with a categorical label—what object is in a photo, what word is spoken in a recording, what topic a document is about. The process of classification discards most of the information in the input and produces a single output (or a probability distribution over values of that single output). The classifier is also often able to ignore many parts of the input. For example, when recognizing an object in a photo, it is usually possible to ignore the background of the photo.

It is possible to ask probabilistic models to do many other tasks. These tasks are often more expensive than classification. Some of them require producing multiple output values. Most require a complete understanding of the entire structure of

---

<sup>1</sup> A *natural image* is an image that might be captured by a camera in a reasonably ordinary environment, as opposed to a synthetically rendered image, a screenshot of a web page, etc.

the input, with no option to ignore sections of it. These tasks include the following:

- Density estimation: given an input  $\mathbf{x}$ , the machine learning system returns an estimate of the true density  $p(\mathbf{x})$  under the data generating distribution. This requires only a single output, but it does require a complete understanding of the entire input. If even one element of the vector is unusual, the system must assign it a low probability.
- Denoising: given a damaged or incorrectly observed input  $\tilde{\mathbf{x}}$ , the machine learning system returns an estimate of the original or correct  $\mathbf{x}$ . For example, the machine learning system might be asked to remove dust or scratches from an old photograph. This requires multiple outputs (every element of the estimated clean example  $\mathbf{x}$ ) and an understanding of the entire input (since even one damaged area will still reveal the final estimate as being damaged).
- Missing value imputation: given the observations of some elements of  $\mathbf{x}$ , the model is asked to return estimates of or a probability distribution over some or all of the unobserved elements of  $\mathbf{x}$ . This requires multiple outputs. Because the model could be asked to restore any of the elements of  $\mathbf{x}$ , it must understand the entire input.
- Sampling: the model generates new samples from the distribution  $p(\mathbf{x})$ . Applications include speech synthesis, i.e. producing new waveforms that sound like natural human speech. This requires multiple output values and a good model of the entire input. If the samples have even one element drawn from the wrong distribution, then the sampling process is wrong.

For an example of a sampling task using small natural images, see Fig. 16.1.

Modeling a rich distribution over thousands or millions of random variables is a challenging task, both computationally and statistically. Suppose we only wanted to model binary variables. This is the simplest possible case, and yet already it seems overwhelming. For a small,  $32 \times 32$  pixel color (RGB) image, there are  $2^{3072}$  possible binary images of this form. This number is over  $10^{800}$  times larger than the estimated number of atoms in the universe.

In general, if we wish to model a distribution over a random vector  $\mathbf{x}$  containing  $n$  discrete variables capable of taking on  $k$  values each, then the naive approach of representing  $P(\mathbf{x})$  by storing a lookup table with one probability value per possible outcome requires  $k^n$  parameters!

This is not feasible for several reasons:

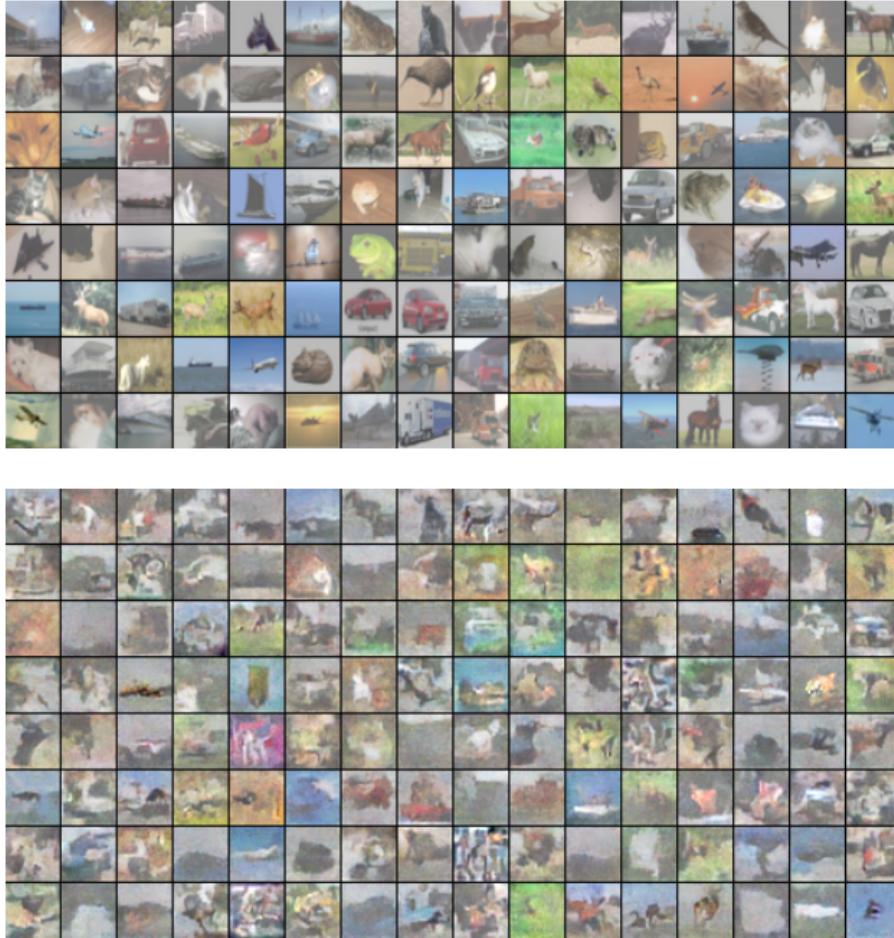


Figure 16.1: Probabilistic modeling of natural images. (*Top*) Example  $32 \times 32$  pixel color images from the CIFAR-10 dataset (Krizhevsky and Hinton, 2009). (*Bottom*) Samples drawn from a structured probabilistic model trained on this dataset. Each sample appears at the same position in the grid as the training example that is closest to it in Euclidean space. This comparison allows us to see that the model is truly synthesizing new images, rather than memorizing the training data. Contrast of both sets of images has been adjusted for display. Figure reproduced with permission from Courville *et al.* (2011).

- **Memory: the cost of storing the representation:** For all but very small values of  $n$  and  $k$ , representing the distribution as a table will require too many values to store.
- **Statistical efficiency:** As the number of parameters in a model increases, so does the amount of training data needed to choose the values of those parameters using a statistical estimator. Because the table-based model has an astronomical number of parameters, it will require an astronomically large training set to fit accurately. Any such model will overfit the training set very badly unless additional assumptions are made linking the different entries in the table (for example, like in back-off or smoothed  $n$ -gram models, Sec. 12.4.1).
- **Runtime: the cost of inference:** Suppose we want to perform an inference task where we use our model of the joint distribution  $P(\mathbf{x})$  to compute some other distribution, such as the marginal distribution  $P(x_1)$  or the conditional distribution  $P(x_2 | x_1)$ . Computing these distributions will require summing across the entire table, so the runtime of these operations is as high as the intractable memory cost of storing the model.
- **Runtime: the cost of sampling:** Likewise, suppose we want to draw a sample from the model. The naive way to do this is to sample some value  $u \sim U(0, 1)$ , then iterate through the table adding up the probability values until they exceed  $u$  and return the outcome whose probability value was added last. This requires reading through the whole table in the worst case, so it has the same exponential cost as the other operations.

The problem with the table-based approach is that we are explicitly modeling every possible kind of interaction between every possible subset of variables. The probability distributions we encounter in real tasks are much simpler than this. Usually, most variables influence each other only indirectly.

For example, consider modeling the finishing times of a team in a relay race. Suppose the team consists of three runners: Alice, Bob and Carol. At the start of the race, Alice carries a baton and begins running around a track. After completing her lap around the track, she hands the baton to Bob. Bob then runs his own lap and hands the baton to Carol, who runs the final lap. We can model each of their finishing times as a continuous random variable. Alice's finishing time does not depend on anyone else's, since she goes first. Bob's finishing time depends on Alice's, because Bob does not have the opportunity to start his lap until Alice has completed hers. If Alice finishes faster, Bob will finish faster, all else being

equal. Finally, Carol’s finishing time depends on both her teammates. If Alice is slow, Bob will probably finish late too. As a consequence, Carol will have quite a late starting time and thus is likely to have a late finishing time as well. However, Carol’s finishing time depends only *indirectly* on Alice’s finishing time via Bob’s. If we already know Bob’s finishing time, we will not be able to estimate Carol’s finishing time better by finding out what Alice’s finishing time was. This means we can model the relay race using only two interactions: Alice’s effect on Bob and Bob’s effect on Carol. We can omit the third, indirect interaction between Alice and Carol from our model.

Structured probabilistic models provide a formal framework for modeling only direct interactions between random variables. This allows the models to have significantly fewer parameters which can in turn be estimated reliably from less data. These smaller models also have dramatically reduced computational cost in terms of storing the model, performing inference in the model, and drawing samples from the model.

## 16.2 Using Graphs to Describe Model Structure

Structured probabilistic models use graphs (in the graph theory sense of “nodes” or “vertices” connected by edges) to represent interactions between random variables. Each node represents a random variable. Each edge represents a direct interaction. These direct interactions imply other, indirect interactions, but only the direct interactions need to be explicitly modeled.

There is more than one way to describe the interactions in a probability distribution using a graph. In the following sections we describe some of the most popular and useful approaches. Graphical models can be largely divided into two categories: models based on directed acyclic graphs, and models based on undirected graphs.

### 16.2.1 Directed Models

One kind of structured probabilistic model is the *directed graphical model*, otherwise known as the *belief network* or *Bayesian network*<sup>2</sup> (Pearl, 1985).

Directed graphical models are called “directed” because their edges are directed,

---

<sup>2</sup> Judea Pearl suggested using the term “Bayesian network” when one wishes to “emphasize the judgmental” nature of the values computed by the network, i.e. to highlight that they usually represent degrees of belief rather than frequencies of events.

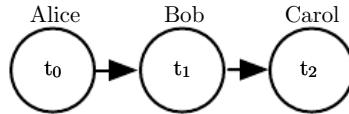


Figure 16.2: A directed graphical model depicting the relay race example. Alice’s finishing time  $t_0$  influences Bob’s finishing time  $t_1$ , because Bob does not get to start running until Alice finishes. Likewise, Carol only gets to start running after Bob finishes, so Bob’s finishing time  $t_1$  directly influences Carol’s finishing time  $t_2$ .

that is, they point from one vertex to another. This direction is represented in the drawing with an arrow. The direction of the arrow indicates which variable’s probability distribution is defined in terms of the other’s. Drawing an arrow from  $a$  to  $b$  means that we define the probability distribution over  $b$  via a conditional distribution, with  $a$  as one of the variables on the right side of the conditioning bar. In other words, the distribution over  $b$  depends on the value of  $a$ .

Continuing with the relay race example from Sec. 16.1, suppose we name Alice’s finishing time  $t_0$ , Bob’s finishing time  $t_1$ , and Carol’s finishing time  $t_2$ . As we saw earlier, our estimate of  $t_1$  depends on  $t_0$ . Our estimate of  $t_2$  depends directly on  $t_1$  but only indirectly on  $t_0$ . We can draw this relationship in a directed graphical model, illustrated in Fig. 16.2.

Formally, a directed graphical model defined on variables  $\mathbf{x}$  is defined by a directed acyclic graph  $\mathcal{G}$  whose vertices are the random variables in the model, and a set of *local conditional probability distributions*  $p(x_i | Pa_{\mathcal{G}}(x_i))$  where  $Pa_{\mathcal{G}}(x_i)$  gives the parents of  $x_i$  in  $\mathcal{G}$ . The probability distribution over  $\mathbf{x}$  is given by

$$p(\mathbf{x}) = \prod_i p(x_i | Pa_{\mathcal{G}}(x_i)). \quad (16.1)$$

In our relay race example, this means that, using the graph drawn in Fig. 16.2,

$$p(t_0, t_1, t_2) = p(t_0)p(t_1 | t_0)p(t_2 | t_1). \quad (16.2)$$

This is our first time seeing a structured probabilistic model in action. We can examine the cost of using it, in order to observe how structured modeling has many advantages relative to unstructured modeling.

Suppose we represented time by discretizing time ranging from minute 0 to minute 10 into 6 second chunks. This would make  $t_0$ ,  $t_1$  and  $t_2$  each be discrete variables with 100 possible values. If we attempted to represent  $p(t_0, t_1, t_2)$  with a table, it would need to store 999,999 values ( $100$  values of  $t_0 \times 100$  values of  $t_1 \times 100$  values of  $t_2$ , minus 1, since the probability of one of the configurations is made

redundant by the constraint that the sum of the probabilities be 1). If instead, we only make a table for each of the conditional probability distributions, then the distribution over  $t_0$  requires 99 values, the table defining  $t_1$  given  $t_0$  requires 9900 values, and so does the table defining  $t_2$  given  $t_1$ . This comes to a total of 19,899 values. This means that using the directed graphical model reduced our number of parameters by a factor of more than 50!

In general, to model  $n$  discrete variables each having  $k$  values, the cost of the single table approach scales like  $O(k^n)$ , as we have observed before. Now suppose we build a directed graphical model over these variables. If  $m$  is the maximum number of variables appearing (on either side of the conditioning bar) in a single conditional probability distribution, then the cost of the tables for the directed model scales like  $O(k^m)$ . As long as we can design a model such that  $m \ll n$ , we get very dramatic savings.

In other words, so long as each variable has few parents in the graph, the distribution can be represented with very few parameters. Some restrictions on the graph structure, such as requiring it to be a tree, can also guarantee that operations like computing marginal or conditional distributions over subsets of variables are efficient.

It is important to realize what kinds of information can and cannot be encoded in the graph. The graph encodes only simplifying assumptions about which variables are conditionally independent from each other. It is also possible to make other kinds of simplifying assumptions. For example, suppose we assume Bob always runs the same regardless of how Alice performed. (In reality, Alice's performance probably influences Bob's performance—depending on Bob's personality, if Alice runs especially fast in a given race, this might encourage Bob to push hard and match her exceptional performance, or it might make him overconfident and lazy). Then the only effect Alice has on Bob's finishing time is that we must add Alice's finishing time to the total amount of time we think Bob needs to run. This observation allows us to define a model with  $O(k)$  parameters instead of  $O(k^2)$ . However, note that  $t_0$  and  $t_1$  are still directly dependent with this assumption, because  $t_1$  represents the absolute time at which Bob finishes, not the total time he himself spends running. This means our graph must still contain an arrow from  $t_0$  to  $t_1$ . The assumption that Bob's personal running time is independent from all other factors cannot be encoded in a graph over  $t_0$ ,  $t_1$ , and  $t_2$ . Instead, we encode this information in the definition of the conditional distribution itself. The conditional distribution is no longer a  $k \times k - 1$  element table indexed by  $t_0$  and  $t_1$  but is now a slightly more complicated formula using only  $k - 1$  parameters. The directed graphical model syntax does not place any constraint on how we define

our conditional distributions. It only defines which variables they are allowed to take in as arguments.

### 16.2.2 Undirected Models

Directed graphical models give us one language for describing structured probabilistic models. Another popular language is that of *undirected models*, otherwise known as *Markov random fields* (MRFs) or *Markov networks* ([Kindermann, 1980](#)). As their name implies, undirected models use graphs whose edges are undirected.

Directed models are most naturally applicable to situations where there is a clear reason to draw each arrow in one particular direction. Often these are situations where we understand the causality and the causality only flows in one direction. One such situation is the relay race example. Earlier runners affect the finishing times of later runners; later runners do not affect the finishing times of earlier runners.

Not all situations we might want to model have such a clear direction to their interactions. When the interactions seem to have no intrinsic direction, or to operate in both directions, it may be more appropriate to use an undirected model.

As an example of such a situation, suppose we want to model a distribution over three binary variables: whether or not you are sick, whether or not your coworker is sick, and whether or not your roommate is sick. As in the relay race example, we can make simplifying assumptions about the kinds of interactions that take place. Assuming that your coworker and your roommate do not know each other, it is very unlikely that one of them will give the other a disease such as a cold directly. This event can be seen as so rare that it is acceptable not to model it. However, it is reasonably likely that either of them could give you a cold, and that you could pass it on to the other. We can model the indirect transmission of a cold from your coworker to your roommate by modeling the transmission of the cold from your coworker to you and the transmission of the cold from you to your roommate.

In this case, it is just as easy for you to cause your roommate to get sick as it is for your roommate to make you sick, so there is not a clean, uni-directional narrative on which to base the model. This motivates using an undirected model. As with directed models, if two nodes in an undirected model are connected by an edge, then the random variables corresponding to those nodes interact with each other directly. Unlike directed models, the edge in an undirected model has no arrow, and is not associated with a conditional probability distribution.

We denote the random variable representing your health as  $h_y$ , the random

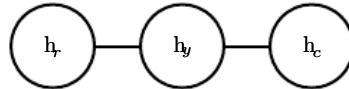


Figure 16.3: An undirected graph representing how your roommate’s health  $h_r$ , your health  $h_y$ , and your work colleague’s health  $h_c$  affect each other. You and your roommate might infect each other with a cold, and you and your work colleague might do the same, but assuming that your roommate and your colleague do not know each other, they can only infect each other indirectly via you.

variable representing your roommate’s health as  $h_r$ , and the random variable representing your colleague’s health as  $h_c$ . See Fig. 16.3 for a drawing of the graph representing this scenario.

Formally, an undirected graphical model is a structured probabilistic model defined on an undirected graph  $\mathcal{G}$ . For each clique  $\mathcal{C}$  in the graph,<sup>3</sup> a factor  $\phi(\mathcal{C})$  (also called a *clique potential*) measures the affinity of the variables in that clique for being in each of their possible joint states. The factors are constrained to be non-negative. Together they define an *unnormalized probability distribution*

$$\tilde{p}(\mathbf{x}) = \prod_{\mathcal{C} \in \mathcal{G}} \phi(\mathcal{C}). \quad (16.3)$$

The unnormalized probability distribution is efficient to work with so long as all the cliques are small. It encodes the idea that states with higher affinity are more likely. However, unlike in a Bayesian network, there is little structure to the definition of the cliques, so there is nothing to guarantee that multiplying them together will yield a valid probability distribution. See Fig. 16.4 for an example of reading factorization information from an undirected graph.

Our example of the cold spreading between you, your roommate, and your colleague contains two cliques. One clique contains  $h_y$  and  $h_c$ . The factor for this clique can be defined by a table, and might have values resembling these:

		$h_y = 0$	$h_y = 1$
		2	1
$h_c = 0$	2		
	1		10

A state of 1 indicates good health, while a state of 0 indicates poor health (having been infected with a cold). Both of you are usually healthy, so the

---

<sup>3</sup>A clique of the graph is a subset of nodes that are all connected to each other by an edge of the graph.

corresponding state has the highest affinity. The state where only one of you is sick has the lowest affinity, because this is a rare state. The state where both of you are sick (because one of you has infected the other) is a higher affinity state, though still not as common as the state where both are healthy.

To complete the model, we would need to also define a similar factor for the clique containing  $h_y$  and  $h_r$ .

### 16.2.3 The Partition Function

While the unnormalized probability distribution is guaranteed to be non-negative everywhere, it is not guaranteed to sum or integrate to 1. To obtain a valid probability distribution, we must use the corresponding normalized probability distribution:<sup>4</sup>

$$p(\mathbf{x}) = \frac{1}{Z} \tilde{p}(\mathbf{x}) \quad (16.4)$$

where  $Z$  is the value that results in the probability distribution summing or integrating to 1:

$$Z = \int \tilde{p}(\mathbf{x}) d\mathbf{x}. \quad (16.5)$$

You can think of  $Z$  as a constant when the  $\phi$  functions are held constant. Note that if the  $\phi$  functions have parameters, then  $Z$  is a function of those parameters. It is common in the literature to write  $Z$  with its arguments omitted to save space. The normalizing constant  $Z$  is known as the *partition function*, a term borrowed from statistical physics.

Since  $Z$  is an integral or sum over all possible joint assignments of the state  $\mathbf{x}$  it is often intractable to compute. In order to be able to obtain the normalized probability distribution of an undirected model, the model structure and the definitions of the  $\phi$  functions must be conducive to computing  $Z$  efficiently. In the context of deep learning,  $Z$  is usually intractable. Due to the intractability of computing  $Z$  exactly, we must resort to approximations. Such approximate algorithms are the topic of Chapter 18.

One important consideration to keep in mind when designing undirected models is that it is possible to specify the factors in such a way that  $Z$  does not exist. This happens if some of the variables in the model are continuous and the integral of  $\tilde{p}$  over their domain diverges. For example, suppose we want to model a single

---

<sup>4</sup>A distribution defined by normalizing a product of clique potentials is also called a *Gibbs distribution*.

scalar variable  $\mathbf{x} \in \mathbb{R}$  with a single clique potential  $\phi(x) = x^2$ . In this case,

$$Z = \int x^2 dx. \quad (16.6)$$

Since this integral diverges, there is no probability distribution corresponding to this choice of  $\phi(x)$ . Sometimes the choice of some parameter of the  $\phi$  functions determines whether the probability distribution is defined. For example, for  $\phi(x; \beta) = \exp(-\beta x^2)$ , the  $\beta$  parameter determines whether  $Z$  exists. Positive  $\beta$  results in a Gaussian distribution over  $\mathbf{x}$  but all other values of  $\beta$  make  $\phi$  impossible to normalize.

One key difference between directed modeling and undirected modeling is that directed models are defined directly in terms of probability distributions from the start, while undirected models are defined more loosely by  $\phi$  functions that are then converted into probability distributions. This changes the intuitions one must develop in order to work with these models. One key idea to keep in mind while working with undirected models is that the domain of each of the variables has dramatic effect on the kind of probability distribution that a given set of  $\phi$  functions corresponds to. For example, consider an  $n$ -dimensional vector-valued random variable  $\mathbf{x}$  and an undirected model parametrized by a vector of biases  $\mathbf{b}$ . Suppose we have one clique for each element of  $\mathbf{x}$ ,  $\phi^{(i)}(\mathbf{x}_i) = \exp(b_i \mathbf{x}_i)$ . What kind of probability distribution does this result in? The answer is that we do not have enough information, because we have not yet specified the domain of  $\mathbf{x}$ . If  $\mathbf{x} \in \mathbb{R}^n$ , then the integral defining  $Z$  diverges and no probability distribution exists. If  $\mathbf{x} \in \{0, 1\}^n$ , then  $p(\mathbf{x})$  factorizes into  $n$  independent distributions, with  $p(\mathbf{x}_i = 1) = \text{sigmoid}(b_i)$ . If the domain of  $\mathbf{x}$  is the set of elementary basis vectors ( $\{[1, 0, \dots, 0], [0, 1, \dots, 0], \dots, [0, 0, \dots, 1]\}$ ) then  $p(\mathbf{x}) = \text{softmax}(\mathbf{b})$ , so a large value of  $b_i$  actually reduces  $p(\mathbf{x}_j = 1)$  for  $j \neq i$ . Often, it is possible to leverage the effect of a carefully chosen domain of a variable in order to obtain complicated behavior from a relatively simple set of  $\phi$  functions. We will explore a practical application of this idea later, in Sec. 20.6.

#### 16.2.4 Energy-Based Models

Many interesting theoretical results about undirected models depend on the assumption that  $\forall \mathbf{x}, \tilde{p}(\mathbf{x}) > 0$ . A convenient way to enforce this condition is to use an *energy-based model* (EBM) where

$$\tilde{p}(\mathbf{x}) = \exp(-E(\mathbf{x})) \quad (16.7)$$

and  $E(\mathbf{x})$  is known as the *energy function*. Because  $\exp(z)$  is positive for all  $z$ , this guarantees that no energy function will result in a probability of zero for any state  $\mathbf{x}$ .

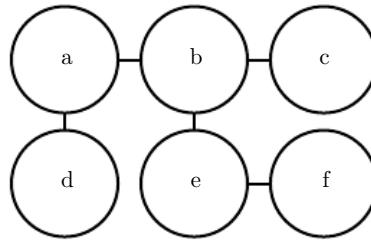


Figure 16.4: This graph implies that  $p(a, b, c, d, e, f)$  can be written as  $\frac{1}{Z} \phi_{a,b}(a, b) \phi_{b,c}(b, c) \phi_{a,d}(a, d) \phi_{b,e}(b, e) \phi_{e,f}(e, f)$  for an appropriate choice of the  $\phi$  functions.

Being completely free to choose the energy function makes learning simpler. If we learned the clique potentials directly, we would need to use constrained optimization to arbitrarily impose some specific minimal probability value. By learning the energy function, we can use unconstrained optimization.<sup>5</sup> The probabilities in an energy-based model can approach arbitrarily close to zero but never reach it.

Any distribution of the form given by Eq. 16.7 is an example of a *Boltzmann distribution*. For this reason, many energy-based models are called *Boltzmann machines* (Fahlman *et al.*, 1983; Ackley *et al.*, 1985; Hinton *et al.*, 1984; Hinton and Sejnowski, 1986). There is no accepted guideline for when to call a model an energy-based model and when to call it a Boltzmann machine. The term Boltzmann machine was first introduced to describe a model with exclusively binary variables, but today many models such as the mean-covariance restricted Boltzmann machine incorporate real-valued variables as well. While Boltzmann machines were originally defined to encompass both models with and without latent variables, the term Boltzmann machine is today most often used to designate models with latent variables, while Boltzmann machines without latent variables are more often called Markov random fields or log-linear models.

Cliques in an undirected graph correspond to factors of the unnormalized probability function. Because  $\exp(a) \exp(b) = \exp(a + b)$ , this means that different cliques in the undirected graph correspond to the different terms of the energy function. In other words, an energy-based model is just a special kind of Markov network: the exponentiation makes each term in the energy function correspond to a factor for a different clique. See Fig. 16.5 for an example of how to read the form of the energy function from an undirected graph structure. One can view an energy-based model with multiple terms in its energy function as being a *product of experts* (Hinton, 1999). Each term in the energy function corresponds to another

<sup>5</sup>For some models, we may still need to use constrained optimization to make sure  $Z$  exists.

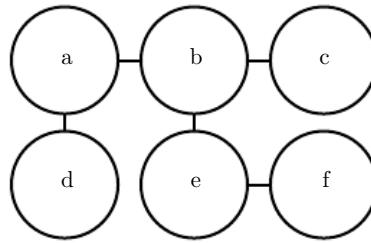


Figure 16.5: This graph implies that  $E(a, b, c, d, e, f)$  can be written as  $E_{a,b}(a, b) + E_{b,c}(b, c) + E_{a,d}(a, d) + E_{b,e}(b, e) + E_{e,f}(e, f)$  for an appropriate choice of the per-clique energy functions. Note that we can obtain the  $\phi$  functions in Fig. 16.4 by setting each  $\phi$  to the exponential of the corresponding negative energy, e.g.,  $\phi_{a,b}(a, b) = \exp(-E(a, b))$ .

factor in the probability distribution. Each term of the energy function can be thought of as an “expert” that determines whether a particular soft constraint is satisfied. Each expert may enforce only one constraint that concerns only a low-dimensional projection of the random variables, but when combined by multiplication of probabilities, the experts together enforce a complicated high-dimensional constraint.

One part of the definition of an energy-based model serves no functional purpose from a machine learning point of view: the  $-$  sign in Eq. 16.7. This  $-$  sign could be incorporated into the definition of  $E$ , or for many functions  $E$  the learning algorithm could simply learn parameters with opposite sign. The  $-$  sign is present primarily to preserve compatibility between the machine learning literature and the physics literature. Many advances in probabilistic modeling were originally developed by statistical physicists, for whom  $E$  refers to actual, physical energy and does not have arbitrary sign. Terminology such as “energy” and “partition function” remains associated with these techniques, even though their mathematical applicability is broader than the physics context in which they were developed. Some machine learning researchers (e.g., Smolensky (1986), who referred to negative energy as *harmony*) have chosen to emit the negation, but this is not the standard convention.

Many algorithms that operate on probabilistic models do not need to compute  $p_{\text{model}}(\mathbf{x})$  but only  $\log \tilde{p}_{\text{model}}(\mathbf{x})$ . For energy-based models with latent variables  $\mathbf{h}$ , these algorithms are sometimes phrased in terms of the negative of this quantity, called the *free energy*:

$$\mathcal{F}(\mathbf{x}) = -\log \sum_{\mathbf{h}} \exp(-E(\mathbf{x}, \mathbf{h})). \quad (16.8)$$

In this book, we usually prefer the more general  $\log \tilde{p}_{\text{model}}(\mathbf{x})$  formulation.

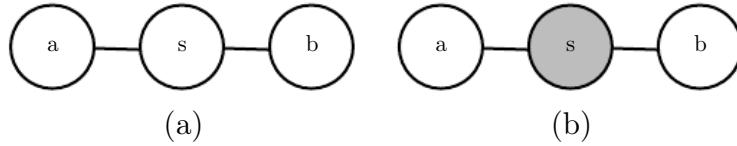


Figure 16.6: (a) The path between random variable  $a$  and random variable  $b$  through  $s$  is active, because  $s$  is not observed. This means that  $a$  and  $b$  are not separated. (b) Here  $s$  is shaded in, to indicate that it is observed. Because the only path between  $a$  and  $b$  is through  $s$ , and that path is inactive, we can conclude that  $a$  and  $b$  are separated given  $s$ .

### 16.2.5 Separation and D-Separation

The edges in a graphical model tell us which variables directly interact. We often need to know which variables *indirectly* interact. Some of these indirect interactions can be enabled or disabled by observing other variables. More formally, we would like to know which subsets of variables are conditionally independent from each other, given the values of other subsets of variables.

Identifying the conditional independences in a graph is very simple in the case of undirected models. In this case, conditional independence implied by the graph is called *separation*. We say that a set of variables  $\mathbb{A}$  is *separated* from another set of variables  $\mathbb{B}$  given a third set of variables  $\mathbb{S}$  if the graph structure implies that  $\mathbb{A}$  is independent from  $\mathbb{B}$  given  $\mathbb{S}$ . If two variables  $a$  and  $b$  are connected by a path involving only unobserved variables, then those variables are not separated. If no path exists between them, or all paths contain an observed variable, then they are separated. We refer to paths involving only unobserved variables as “active” and paths including an observed variable as “inactive.”

When we draw a graph, we can indicate observed variables by shading them in. See Fig. 16.6 for a depiction of how active and inactive paths in an undirected model look when drawn in this way. See Fig. 16.7 for an example of reading separation from an undirected graph.

Similar concepts apply to directed models, except that in the context of directed models, these concepts are referred to as *d-separation*. The “d” stands for “dependence.” D-separation for directed graphs is defined the same as separation for undirected graphs: We say that a set of variables  $\mathbb{A}$  is d-separated from another set of variables  $\mathbb{B}$  given a third set of variables  $\mathbb{S}$  if the graph structure implies that  $\mathbb{A}$  is independent from  $\mathbb{B}$  given  $\mathbb{S}$ .

As with undirected models, we can examine the independences implied by the graph by looking at what active paths exist in the graph. As before, two variables are dependent if there is an active path between them, and d-separated if no such

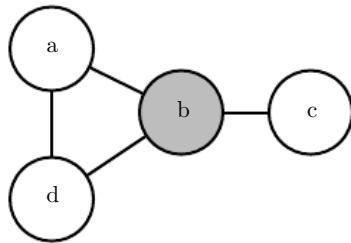


Figure 16.7: An example of reading separation properties from an undirected graph. Here  $b$  is shaded to indicate that it is observed. Because observing  $b$  blocks the only path from  $a$  to  $c$ , we say that  $a$  and  $c$  are separated from each other given  $b$ . The observation of  $b$  also blocks one path between  $a$  and  $d$ , but there is a second, active path between them. Therefore,  $a$  and  $d$  are not separated given  $b$ .

path exists. In directed nets, determining whether a path is active is somewhat more complicated. See Fig. 16.8 for a guide to identifying active paths in a directed model. See Fig. 16.9 for an example of reading some properties from a graph.

It is important to remember that separation and  $d$ -separation tell us only about those conditional independences **that are implied by the graph**. There is no requirement that the graph imply all independences that are present. In particular, it is always legitimate to use the complete graph (the graph with all possible edges) to represent any distribution. In fact, some distributions contain independences that are not possible to represent with existing graphical notation. *Context-specific independences* are independences that are present dependent on the value of some variables in the network. For example, consider a model of three binary variables:  $a$ ,  $b$  and  $c$ . Suppose that when  $a$  is 0,  $b$  and  $c$  are independent, but when  $a$  is 1,  $b$  is deterministically equal to  $c$ . Encoding the behavior when  $a = 1$  requires an edge connecting  $b$  and  $c$ . The graph then fails to indicate that  $b$  and  $c$  are independent when  $a = 0$ .

In general, a graph will never imply that an independence exists when it does not. However, a graph may fail to encode an independence.

### 16.2.6 Converting between Undirected and Directed Graphs

We often refer to a specific machine learning model as being undirected or directed. For example, we typically refer to RBMs as undirected and sparse coding as directed. This choice of wording can be somewhat misleading, because no probabilistic model is inherently directed or undirected. Instead, some models are most easily **described** using a directed graph, or most easily described using an undirected graph.

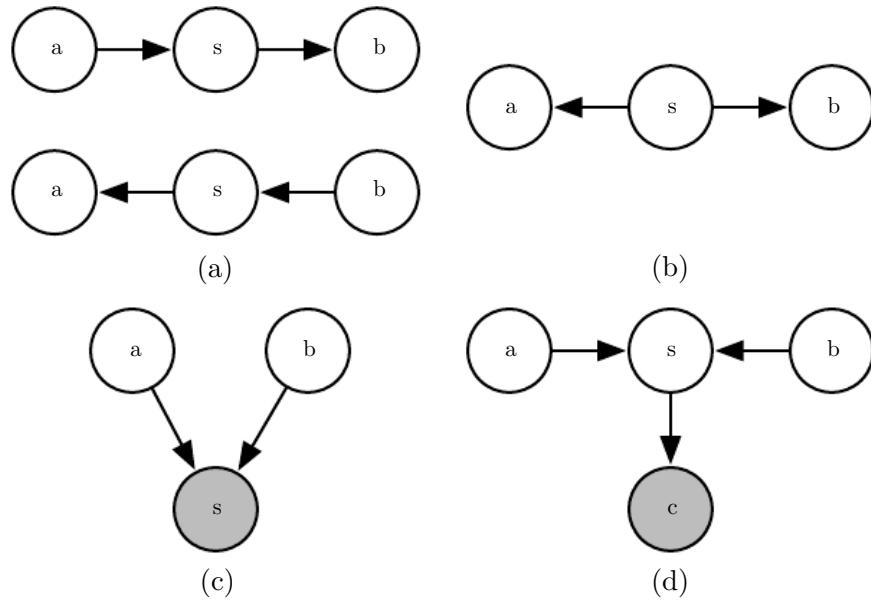


Figure 16.8: All of the kinds of active paths of length two that can exist between random variables  $a$  and  $b$ . (a) Any path with arrows proceeding directly from  $a$  to  $b$  or vice versa. This kind of path becomes blocked if  $s$  is observed. We have already seen this kind of path in the relay race example. (b)  $a$  and  $b$  are connected by a *common cause*  $s$ . For example, suppose  $s$  is a variable indicating whether or not there is a hurricane and  $a$  and  $b$  measure the wind speed at two different nearby weather monitoring outposts. If we observe very high winds at station  $a$ , we might expect to also see high winds at  $b$ . This kind of path can be blocked by observing  $s$ . If we already know there is a hurricane, we expect to see high winds at  $b$ , regardless of what is observed at  $a$ . A lower than expected wind at  $a$  (for a hurricane) would not change our expectation of winds at  $b$  (knowing there is a hurricane). However, if  $s$  is not observed, then  $a$  and  $b$  are dependent, i.e., the path is active. (c)  $a$  and  $b$  are both parents of  $s$ . This is called a *V-structure* or *the collider case*. The V-structure causes  $a$  and  $b$  to be related by the *explaining away effect*. In this case, the path is actually active when  $s$  is observed. For example, suppose  $s$  is a variable indicating that your colleague is not at work. The variable  $a$  represents her being sick, while  $b$  represents her being on vacation. If you observe that she is not at work, you can presume she is probably sick or on vacation, but it is not especially likely that both have happened at the same time. If you find out that she is on vacation, this fact is sufficient to *explain* her absence. You can infer that she is probably not also sick. (d) The explaining away effect happens even if any descendant of  $s$  is observed! For example, suppose that  $c$  is a variable representing whether you have received a report from your colleague. If you notice that you have not received the report, this increases your estimate of the probability that she is not at work today, which in turn makes it more likely that she is either sick or on vacation. The only way to block a path through a V-structure is to observe none of the descendants of the shared child.

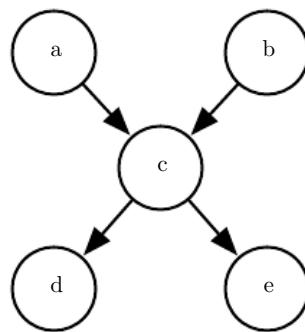


Figure 16.9: From this graph, we can read out several d-separation properties. Examples include:

- a and b are d-separated given the empty set.
- a and e are d-separated given c.
- d and e are d-separated given c.

We can also see that some variables are no longer d-separated when we observe some variables:

- a and b are not d-separated given c.
- a and b are not d-separated given d.

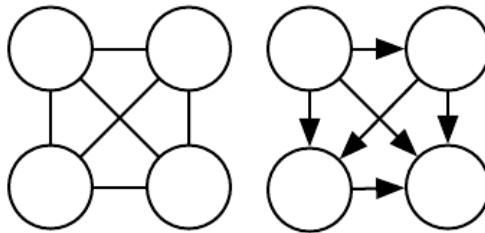


Figure 16.10: Examples of complete graphs, which can describe any probability distribution. Here we show examples with four random variables. (*Left*) The complete undirected graph. In the undirected case, the complete graph is unique. (*Right*) A complete directed graph. In the directed case, there is not a unique complete graph. We choose an ordering of the variables and draw an arc from each variable to every variable that comes after it in the ordering. There are thus a factorial number of complete graphs for every set of random variables. In this example we order the variables from left to right, top to bottom.

Directed models and undirected models both have their advantages and disadvantages. Neither approach is clearly superior and universally preferred. Instead, we should choose which language to use for each task. This choice will partially depend on which probability distribution we wish to describe. We may choose to use either directed modeling or undirected modeling based on which approach can capture the most independences in the probability distribution or which approach uses the fewest edges to describe the distribution. There are other factors that can affect the decision of which language to use. Even while working with a single probability distribution, we may sometimes switch between different modeling languages. Sometimes a different language becomes more appropriate if we observe a certain subset of variables, or if we wish to perform a different computational task. For example, the directed model description often provides a straightforward approach to efficiently draw samples from the model (described in Sec. 16.3) while the undirected model formulation is often useful for deriving approximate inference procedures (as we will see in Chapter 19, where the role of undirected models is highlighted in Eq. 19.56).

Every probability distribution can be represented by either a directed model or by an undirected model. In the worst case, one can always represent any distribution by using a “complete graph.” In the case of a directed model, the complete graph is any directed acyclic graph where we impose some ordering on the random variables, and each variable has all other variables that precede it in the ordering as its ancestors in the graph. For an undirected model, the complete graph is simply a graph containing a single clique encompassing all of the variables. See Fig. 16.10 for an example.

Of course, the utility of a graphical model is that the graph implies that some variables do not interact directly. The complete graph is not very useful because it does not imply any independences.

When we represent a probability distribution with a graph, we want to choose a graph that implies as many independences as possible, without implying any independences that do not actually exist.

From this point of view, some distributions can be represented more efficiently using directed models, while other distributions can be represented more efficiently using undirected models. In other words, directed models can encode some independences that undirected models cannot encode, and vice versa.

Directed models are able to use one specific kind of substructure that undirected models cannot represent perfectly. This substructure is called an *immorality*. The structure occurs when two random variables  $a$  and  $b$  are both parents of a third random variable  $c$ , and there is no edge directly connecting  $a$  and  $b$  in either direction. (The name “immorality” may seem strange; it was coined in the graphical models literature as a joke about unmarried parents.) To convert a directed model with graph  $\mathcal{D}$  into an undirected model, we need to create a new graph  $\mathcal{U}$ . For every pair of variables  $x$  and  $y$ , we add an undirected edge connecting  $x$  and  $y$  to  $\mathcal{U}$  if there is a directed edge (in either direction) connecting  $x$  and  $y$  in  $\mathcal{D}$  or if  $x$  and  $y$  are both parents in  $\mathcal{D}$  of a third variable  $z$ . The resulting  $\mathcal{U}$  is known as a *moralized graph*. See Fig. 16.11 for examples of converting directed models to undirected models via moralization.

Likewise, undirected models can include substructures that no directed model can represent perfectly. Specifically, a directed graph  $\mathcal{D}$  cannot capture all of the conditional independences implied by an undirected graph  $\mathcal{U}$  if  $\mathcal{U}$  contains a *loop* of length greater than three, unless that loop also contains a *chord*. A loop is a sequence of variables connected by undirected edges, with the last variable in the sequence connected back to the first variable in the sequence. A chord is a connection between any two non-consecutive variables in the sequence defining a loop. If  $\mathcal{U}$  has loops of length four or greater and does not have chords for these loops, we must add the chords before we can convert it to a directed model. Adding these chords discards some of the independence information that was encoded in  $\mathcal{U}$ . The graph formed by adding chords to  $\mathcal{U}$  is known as a *chordal* or *triangulated* graph, because all the loops can now be described in terms of smaller, triangular loops. To build a directed graph  $\mathcal{D}$  from the chordal graph, we need to also assign directions to the edges. When doing so, we must not create a directed cycle in  $\mathcal{D}$ , or the result does not define a valid directed probabilistic model. One way to assign directions to the edges in  $\mathcal{D}$  is to impose an ordering on the random

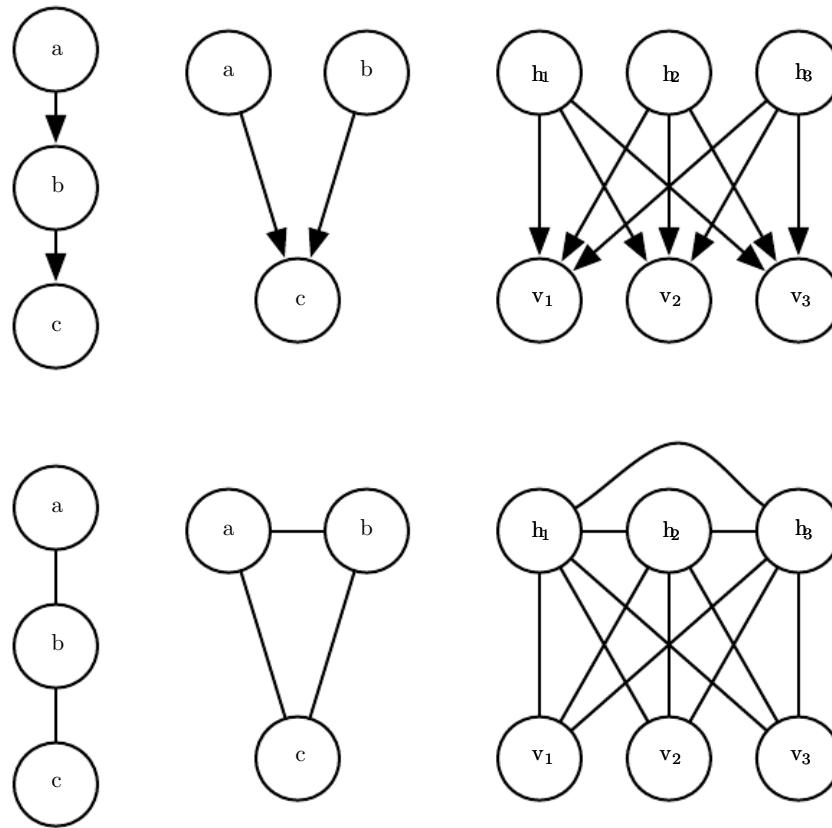


Figure 16.11: Examples of converting directed models (top row) to undirected models (bottom row) by constructing moralized graphs. *(Left)* This simple chain can be converted to a moralized graph merely by replacing its directed edges with undirected edges. The resulting undirected model implies exactly the same set of independences and conditional independences. *(Center)* This graph is the simplest directed model that cannot be converted to an undirected model without losing some independences. This graph consists entirely of a single immorality. Because a and b are parents of c, they are connected by an active path when c is observed. To capture this dependence, the undirected model must include a clique encompassing all three variables. This clique fails to encode the fact that  $a \perp b$ . *(Right)* In general, moralization may add many edges to the graph, thus losing many implied independences. For example, this sparse coding graph requires adding moralizing edges between every pair of hidden units, thus introducing a quadratic number of new direct dependences.

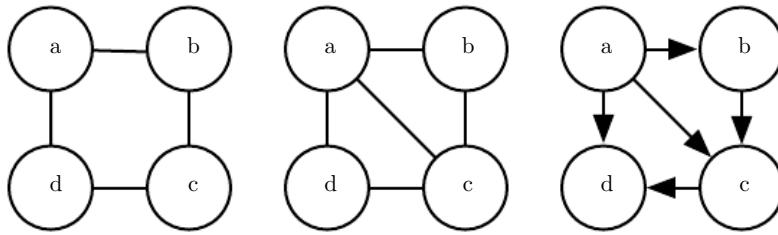


Figure 16.12: Converting an undirected model to a directed model. (*Left*) This undirected model cannot be converted directly to a directed model because it has a loop of length four with no chords. Specifically, the undirected model encodes two different independences that no directed model can capture simultaneously:  $a \perp\!\!\!\perp c \mid \{b, d\}$  and  $b \perp\!\!\!\perp d \mid \{a, c\}$ . (*Center*) To convert the undirected model to a directed model, we must triangulate the graph, by ensuring that all loops of greater than length three have a chord. To do so, we can either add an edge connecting  $a$  and  $c$  or we can add an edge connecting  $b$  and  $d$ . In this example, we choose to add the edge connecting  $a$  and  $c$ . (*Right*) To finish the conversion process, we must assign a direction to each edge. When doing so, we must not create any directed cycles. One way to avoid directed cycles is to impose an ordering over the nodes, and always point each edge from the node that comes earlier in the ordering to the node that comes later in the ordering. In this example, we use the variable names to impose alphabetical order.

variables, then point each edge from the node that comes earlier in the ordering to the node that comes later in the ordering. See Fig. 16.12 for a demonstration.

### 16.2.7 Factor Graphs

*Factor graphs* are another way of drawing undirected models that resolve an ambiguity in the graphical representation of standard undirected model syntax. In an undirected model, the scope of every  $\phi$  function must be a subset of some clique in the graph. However, it is not necessary that there exist any  $\phi$  whose scope contains the entirety of every clique. Factor graphs explicitly represent the scope of each  $\phi$  function. Specifically, a factor graph is a graphical representation of an undirected model that consists of a bipartite undirected graph. Some of the nodes are drawn as circles. These nodes correspond to random variables as in a standard undirected model. The rest of the nodes are drawn as squares. These nodes correspond to the factors  $\phi$  of the unnormalized probability distribution. Variables and factors may be connected with undirected edges. A variable and a factor are connected in the graph if and only if the variable is one of the arguments to the factor in the unnormalized probability distribution. No factor may be connected to another factor in the graph, nor can a variable be connected to a

variable. See Fig. 16.13 for an example of how factor graphs can resolve ambiguity in the interpretation of undirected networks.

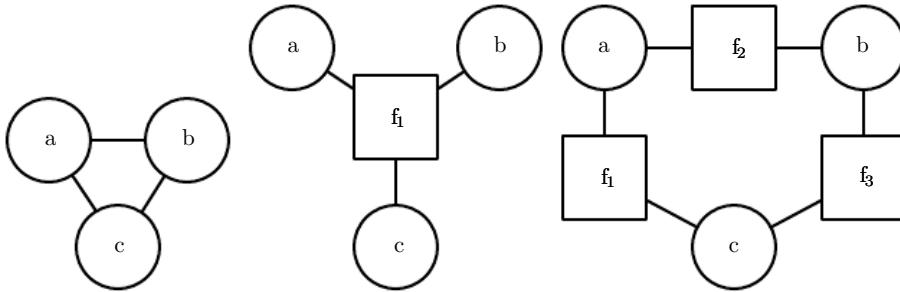


Figure 16.13: An example of how a factor graph can resolve ambiguity in the interpretation of undirected networks. (*Left*) An undirected network with a clique involving three variables: a, b and c. (*Center*) A factor graph corresponding to the same undirected model. This factor graph has one factor over all three variables. (*Right*) Another valid factor graph for the same undirected model. This factor graph has three factors, each over only two variables. Representation, inference, and learning are all asymptotically cheaper in (c) compared to (b), even though both require the same undirected graph to represent.

### 16.3 Sampling from Graphical Models

Graphical models also facilitate the task of drawing samples from a model.

One advantage of directed graphical models is that a simple and efficient procedure called *ancestral sampling* can produce a sample from the joint distribution represented by the model.

The basic idea is to sort the variables  $x_i$  in the graph into a topological ordering, so that for all  $i$  and  $j$ ,  $j$  is greater than  $i$  if  $x_i$  is a parent of  $x_j$ . The variables can then be sampled in this order. In other words, we first sample  $x_1 \sim P(x_1)$ , then sample  $P(x_2 | Pa_G(x_2))$ , and so on, until finally we sample  $P(x_n | Pa_G(x_n))$ . So long as each conditional distribution  $p(x_i | Pa_G(x_i))$  is easy to sample from, then the whole model is easy to sample from. The topological sorting operation guarantees that we can read the conditional distributions in Eq. 16.1 and sample from them in order. Without the topological sorting, we might attempt to sample a variable before its parents are available.

For some graphs, more than one topological ordering is possible. Ancestral sampling may be used with any of these topological orderings.

Ancestral sampling is generally very fast (assuming sampling from each conditional is easy) and convenient.

One drawback to ancestral sampling is that it only applies to directed graphical models. Another drawback is that it does not support every conditional sampling operation. When we wish to sample from a subset of the variables in a directed graphical model, given some other variables, we often require that all the conditioning variables come earlier than the variables to be sampled in the ordered graph. In this case, we can sample from the local conditional probability distributions specified by the model distribution. Otherwise, the conditional distributions we need to sample from are the posterior distributions given the observed variables. These posterior distributions are usually not explicitly specified and parametrized in the model. Inferring these posterior distributions can be costly. In models where this is the case, ancestral sampling is no longer efficient.

Unfortunately, ancestral sampling is only applicable to directed models. We can sample from undirected models by converting them to directed models, but this often requires solving intractable inference problems (to determine the marginal distribution over the root nodes of the new directed graph) or requires introducing so many edges that the resulting directed model becomes intractable. Sampling from an undirected model without first converting it to a directed model seems to require resolving cyclical dependencies. Every variable interacts with every other variable, so there is no clear beginning point for the sampling process. Unfortunately, drawing samples from an undirected graphical model is an expensive, multi-pass process. The conceptually simplest approach is *Gibbs sampling*. Suppose we have a graphical model over an  $n$ -dimensional vector of random variables  $\mathbf{x}$ . We iteratively visit each variable  $x_i$  and draw a sample conditioned on all of the other variables, from  $p(x_i \mid \mathbf{x}_{-i})$ . Due to the separation properties of the graphical model, we can equivalently condition on only the neighbors of  $x_i$ . Unfortunately, after we have made one pass through the graphical model and sampled all  $n$  variables, we still do not have a fair sample from  $p(\mathbf{x})$ . Instead, we must repeat the process and resample all  $n$  variables using the updated values of their neighbors. Asymptotically, after many repetitions, this process converges to sampling from the correct distribution. It can be difficult to determine when the samples have reached a sufficiently accurate approximation of the desired distribution. Sampling techniques for undirected models are an advanced topic, covered in more detail in Chapter 17.

## 16.4 Advantages of Structured Modeling

The primary advantage of using structured probabilistic models is that they allow us to dramatically reduce the cost of representing probability distributions as well

as learning and inference. Sampling is also accelerated in the case of directed models, while the situation can be complicated with undirected models. The primary mechanism that allows all of these operations to use less runtime and memory is choosing to not model certain interactions. Graphical models convey information by leaving edges out. Anywhere there is not an edge, the model specifies the assumption that we do not need to model a direct interaction.

A less quantifiable benefit of using structured probabilistic models is that they allow us to explicitly separate representation of knowledge from learning of knowledge or inference given existing knowledge. This makes our models easier to develop and debug. We can design, analyze, and evaluate learning algorithms and inference algorithms that are applicable to broad classes of graphs. Independently, we can design models that capture the relationships we believe are important in our data. We can then combine these different algorithms and structures and obtain a Cartesian product of different possibilities. It would be much more difficult to design end-to-end algorithms for every possible situation.

## 16.5 Learning about Dependencies

A good generative model needs to accurately capture the distribution over the observed or “visible” variables  $\mathbf{v}$ . Often the different elements of  $\mathbf{v}$  are highly dependent on each other. In the context of deep learning, the approach most commonly used to model these dependencies is to introduce several latent or “hidden” variables,  $\mathbf{h}$ . The model can then capture dependencies between any pair of variables  $v_i$  and  $v_j$  indirectly, via direct dependencies between  $v_i$  and  $\mathbf{h}$ , and direct dependencies between  $\mathbf{h}$  and  $v_j$ .

A good model of  $\mathbf{v}$  which did not contain any latent variables would need to have very large numbers of parents per node in a Bayesian network or very large cliques in a Markov network. Just representing these higher order interactions is costly—both in a computational sense, because the number of parameters that must be stored in memory scales exponentially with the number of members in a clique, but also in a statistical sense, because this exponential number of parameters requires a wealth of data to estimate accurately.

When the model is intended to capture dependencies between visible variables with direct connections, it is usually infeasible to connect all variables, so the graph must be designed to connect those variables that are tightly coupled and omit edges between other variables. An entire field of machine learning called *structure learning* is devoted to this problem. For a good reference on structure learning, see (Koller and Friedman, 2009). Most structure learning techniques are a form of

greedy search. A structure is proposed, a model with that structure is trained, then given a score. The score rewards high training set accuracy and penalizes model complexity. Candidate structures with a small number of edges added or removed are then proposed as the next step of the search. The search proceeds to a new structure that is expected to increase the score.

Using latent variables instead of adaptive structure avoids the need to perform discrete searches and multiple rounds of training. A fixed structure over visible and hidden variables can use direct interactions between visible and hidden units to impose indirect interactions between visible units. Using simple parameter learning techniques we can learn a model with a fixed structure that imputes the right structure on the marginal  $p(\mathbf{v})$ .

Latent variables have advantages beyond their role in efficiently capturing  $p(\mathbf{v})$ . The new variables  $\mathbf{h}$  also provide an alternative representation for  $\mathbf{v}$ . For example, as discussed in Sec. 3.9.6, the mixture of Gaussians model learns a latent variable that corresponds to which category of examples the input was drawn from. This means that the latent variable in a mixture of Gaussians model can be used to do classification. In Chapter 14 we saw how simple probabilistic models like sparse coding learn latent variables that can be used as input features for a classifier, or as coordinates along a manifold. Other models can be used in this same way, but deeper models and models with different kinds of interactions can create even richer descriptions of the input. Many approaches accomplish feature learning by learning latent variables. Often, given some model of  $\mathbf{v}$  and  $\mathbf{h}$ , experimental observations show that  $\mathbb{E}[\mathbf{h} | \mathbf{v}]$  or  $\text{argmax}_{\mathbf{h}} p(\mathbf{h}, \mathbf{v})$  is a good feature mapping for  $\mathbf{v}$ .

## 16.6 Inference and Approximate Inference

One of the main ways we can use a probabilistic model is to ask questions about how variables are related to each other. Given a set of medical tests, we can ask what disease a patient might have. In a latent variable model, we might want to extract features  $\mathbb{E}[\mathbf{h} | \mathbf{v}]$  describing the observed variables  $\mathbf{v}$ . Sometimes we need to solve such problems in order to perform other tasks. We often train our models using the principle of maximum likelihood. Because

$$\log p(\mathbf{v}) = \mathbb{E}_{\mathbf{h} \sim p(\mathbf{h}|\mathbf{v})} [\log p(\mathbf{h}, \mathbf{v}) - \log p(\mathbf{h} | \mathbf{v})], \quad (16.9)$$

we often want to compute  $p(\mathbf{h} | \mathbf{v})$  in order to implement a learning rule. All of these are examples of *inference* problems in which we must predict the value of

some variables given other variables, or predict the probability distribution over some variables given the value of other variables.

Unfortunately, for most interesting deep models, these inference problems are intractable, even when we use a structured graphical model to simplify them. The graph structure allows us to represent complicated, high-dimensional distributions with a reasonable number of parameters, but the graphs used for deep learning are usually not restrictive enough to also allow efficient inference.

It is straightforward to see that computing the marginal probability of a general graphical model is  $\#P$  hard. The complexity class  $\#P$  is a generalization of the complexity class NP. Problems in NP require determining only whether a problem has a solution and finding a solution if one exists. Problems in  $\#P$  require counting the number of solutions. To construct a worst-case graphical model, imagine that we define a graphical model over the binary variables in a 3-SAT problem. We can impose a uniform distribution over these variables. We can then add one binary latent variable per clause that indicates whether each clause is satisfied. We can then add another latent variable indicating whether all of the clauses are satisfied. This can be done without making a large clique, by building a reduction tree of latent variables, with each node in the tree reporting whether two other variables are satisfied. The leaves of this tree are the variables for each clause. The root of the tree reports whether the entire problem is satisfied. Due to the uniform distribution over the literals, the marginal distribution over the root of the reduction tree specifies what fraction of assignments satisfy the problem. While this is a contrived worst-case example, NP hard graphs commonly arise in practical real-world scenarios.

This motivates the use of approximate inference. In the context of deep learning, this usually refers to variational inference, in which we approximate the true distribution  $p(\mathbf{h} | \mathbf{v})$  by seeking an approximate distribution  $q(\mathbf{h}|\mathbf{v})$  that is as close to the true one as possible. This and other techniques are described in depth in Chapter 19.

## 16.7 The Deep Learning Approach to Structured Probabilistic Models

Deep learning practitioners generally use the same basic computational tools as other machine learning practitioners who work with structured probabilistic models. However, in the context of deep learning, we usually make different design decisions about how to combine these tools, resulting in overall algorithms and models that

have a very different flavor from more traditional graphical models.

Deep learning does not always involve especially deep graphical models. In the context of graphical models, we can define the depth of a model in terms of the graphical model graph rather than the computational graph. We can think of a latent variable  $h_i$  as being at depth  $j$  if the shortest path from  $h_i$  to an observed variable is  $j$  steps. We usually describe the depth of the model as being the greatest depth of any such  $h_i$ . This kind of depth is different from the depth induced by the computational graph. Many generative models used for deep learning have no latent variables or only one layer of latent variables, but use deep computational graphs to define the conditional distributions within a model.

Deep learning essentially always makes use of the idea of distributed representations. Even shallow models used for deep learning purposes (such as pretraining shallow models that will later be composed to form deep ones) nearly always have a single, large layer of latent variables. Deep learning models typically have more latent variables than observed variables. Complicated nonlinear interactions between variables are accomplished via indirect connections that flow through multiple latent variables.

By contrast, traditional graphical models usually contain mostly variables that are at least occasionally observed, even if many of the variables are missing at random from some training examples. Traditional models mostly use higher-order terms and structure learning to capture complicated nonlinear interactions between variables. If there are latent variables, they are usually few in number.

The way that latent variables are designed also differs in deep learning. The deep learning practitioner typically does not intend for the latent variables to take on any specific semantics ahead of time—the training algorithm is free to invent the concepts it needs to model a particular dataset. The latent variables are usually not very easy for a human to interpret after the fact, though visualization techniques may allow some rough characterization of what they represent. When latent variables are used in the context of traditional graphical models, they are often designed with some specific semantics in mind—the topic of a document, the intelligence of a student, the disease causing a patient’s symptoms, etc. These models are often much more interpretable by human practitioners and often have more theoretical guarantees, yet are less able to scale to complex problems and are not reusable in as many different contexts as deep models.

Another obvious difference is the kind of connectivity typically used in the deep learning approach. Deep graphical models typically have large groups of units that are all connected to other groups of units, so that the interactions between two groups may be described by a single matrix. Traditional graphical models

have very few connections and the choice of connections for each variable may be individually designed. The design of the model structure is tightly linked with the choice of inference algorithm. Traditional approaches to graphical models typically aim to maintain the tractability of exact inference. When this constraint is too limiting, a popular approximate inference algorithm is an algorithm called *loopy belief propagation*. Both of these approaches often work well with very sparsely connected graphs. By comparison, models used in deep learning tend to connect each visible unit  $v_i$  to very many hidden units  $h_j$ , so that  $\mathbf{h}$  can provide a distributed representation of  $v_i$  (and probably several other observed variables too). Distributed representations have many advantages, but from the point of view of graphical models and computational complexity, distributed representations have the disadvantage of usually yielding graphs that are not sparse enough for the traditional techniques of exact inference and loopy belief propagation to be relevant. As a consequence, one of the most striking differences between the larger graphical models community and the deep graphical models community is that loopy belief propagation is almost never used for deep learning. Most deep models are instead designed to make Gibbs sampling or variational inference algorithms efficient. Another consideration is that deep learning models contain a very large number of latent variables, making efficient numerical code essential. This provides an additional motivation, besides the choice of high-level inference algorithm, for grouping the units into layers with a matrix describing the interaction between two layers. This allows the individual steps of the algorithm to be implemented with efficient matrix product operations, or sparsely connected generalizations, like block diagonal matrix products or convolutions.

Finally, the deep learning approach to graphical modeling is characterized by a marked tolerance of the unknown. Rather than simplifying the model until all quantities we might want can be computed exactly, we increase the power of the model until it is just barely possible to train or use. We often use models whose marginal distributions cannot be computed, and are satisfied simply to draw approximate samples from these models. We often train models with an intractable objective function that we cannot even approximate in a reasonable amount of time, but we are still able to approximately train the model if we can efficiently obtain an estimate of the gradient of such a function. The deep learning approach is often to figure out what the minimum amount of information we absolutely need is, and then to figure out how to get a reasonable approximation of that information as quickly as possible.

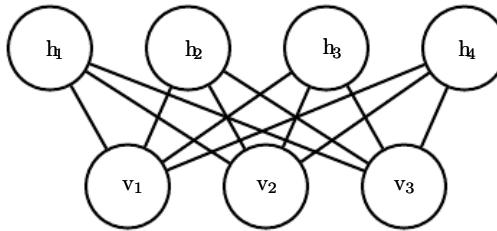


Figure 16.14: An RBM drawn as a Markov network.

### 16.7.1 Example: The Restricted Boltzmann Machine

The *restricted Boltzmann machine* (RBM) (Smolensky, 1986) or *harmonium* is the quintessential example of how graphical models are used for deep learning. The RBM is not itself a deep model. Instead, it has a single layer of latent variables that may be used to learn a representation for the input. In Chapter 20, we will see how RBMs can be used to build many deeper models. Here, we show how the RBM exemplifies many of the practices used in a wide variety of deep graphical models: its units are organized into large groups called layers, the connectivity between layers is described by a matrix, the connectivity is relatively dense, the model is designed to allow efficient Gibbs sampling, and the emphasis of the model design is on freeing the training algorithm to learn latent variables whose semantics were not specified by the designer. Later, in Sec. 20.2, we will revisit the RBM in more detail.

The canonical RBM is an energy-based model with binary visible and hidden units. Its energy function is

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{b}^\top \mathbf{v} - \mathbf{c}^\top \mathbf{h} - \mathbf{v}^\top \mathbf{W} \mathbf{h}, \quad (16.10)$$

where  $\mathbf{b}$ ,  $\mathbf{c}$ , and  $\mathbf{W}$  are unconstrained, real-valued, learnable parameters. We can see that the model is divided into two groups of units:  $\mathbf{v}$  and  $\mathbf{h}$ , and the interaction between them is described by a matrix  $\mathbf{W}$ . The model is depicted graphically in Fig. 16.14. As this figure makes clear, an important aspect of this model is that there are no direct interactions between any two visible units or between any two hidden units (hence the “restricted,” a general Boltzmann machine may have arbitrary connections).

The restrictions on the RBM structure yield the nice properties

$$p(\mathbf{h} \mid \mathbf{v}) = \prod_i p(h_i \mid \mathbf{v}) \quad (16.11)$$

and

$$p(\mathbf{v} \mid \mathbf{h}) = \prod_i p(v_i \mid \mathbf{h}). \quad (16.12)$$



Figure 16.15: Samples from a trained RBM, and its weights. Image reproduced with permission from LISA (2008). (*Left*) Samples from a model trained on MNIST, drawn using Gibbs sampling. Each column is a separate Gibbs sampling process. Each row represents the output of another 1,000 steps of Gibbs sampling. Successive samples are highly correlated with one another. (*Right*) The corresponding weight vectors. Compare this to the samples and weights of a linear factor model, shown in Fig. 13.2. The samples here are much better because the RBM prior  $p(\mathbf{h})$  is not constrained to be factorial. The RBM can learn which features should appear together when sampling. On the other hand, the RBM posterior  $p(\mathbf{h} \mid \mathbf{v})$  is factorial, while the sparse coding posterior  $p(\mathbf{h} \mid \mathbf{v})$  is not, so the sparse coding model may be better for feature extraction. Other models are able to have both a non-factorial  $p(\mathbf{h})$  and a non-factorial  $p(\mathbf{h} \mid \mathbf{v})$ .

The individual conditionals are simple to compute as well. For the binary RBM we obtain:

$$P(h_i = 1 \mid \mathbf{v}) = \sigma \left( \mathbf{v}^\top \mathbf{W}_{:,i} + b_i \right), \quad (16.13)$$

$$P(h_i = 0 \mid \mathbf{v}) = 1 - \sigma \left( \mathbf{v}^\top \mathbf{W}_{:,i} + b_i \right). \quad (16.14)$$

Together these properties allow for efficient *block Gibbs* sampling, which alternates between sampling all of  $\mathbf{h}$  simultaneously and sampling all of  $\mathbf{v}$  simultaneously. Samples generated by Gibbs sampling from an RBM model are shown in Fig. 16.15.

Since the energy function itself is just a linear function of the parameters, it is easy to take derivatives of the energy function. For example,

$$\frac{\partial}{\partial W_{i,j}} E(\mathbf{v}, \mathbf{h}) = -v_i h_j. \quad (16.15)$$

These two properties—efficient Gibbs sampling and efficient derivatives—make training convenient. In Chapter 18, we will see that undirected models may be trained by computing such derivatives applied to samples from the model.

Training the model induces a representation  $\mathbf{h}$  of the data  $\mathbf{v}$ . We can often use  $\mathbb{E}_{\mathbf{h} \sim p(\mathbf{h} \mid \mathbf{v})} [\mathbf{h}]$  as a set of features to describe  $\mathbf{v}$ .

Overall, the RBM demonstrates the typical deep learning approach to graphical models: representation learning accomplished via layers of latent variables, combined with efficient interactions between layers parametrized by matrices.

The language of graphical models provides an elegant, flexible and clear language for describing probabilistic models. In the chapters ahead, we use this language, among other perspectives, to describe a wide variety of deep probabilistic models.

# Chapter 17

## Monte Carlo Methods

Randomized algorithms fall into two rough categories: Las Vegas algorithms and Monte Carlo algorithms. Las Vegas algorithms always return precisely the correct answer (or report that they failed). These algorithms consume a random amount of resources, usually memory or time. In contrast, Monte Carlo algorithms return answers with a random amount of error. The amount of error can typically be reduced by expending more resources (usually running time and memory). For any fixed computational budget, a Monte Carlo algorithm can provide an approximate answer.

Many problems in machine learning are so difficult that we can never expect to obtain precise answers to them. This excludes precise deterministic algorithms and Las Vegas algorithms. Instead, we must use deterministic approximate algorithms or Monte Carlo approximations. Both approaches are ubiquitous in machine learning. In this chapter, we focus on Monte Carlo methods.

### 17.1 Sampling and Monte Carlo Methods

Many important technologies used to accomplish machine learning goals are based on drawing samples from some probability distribution and using these samples to form a Monte Carlo estimate of some desired quantity.

#### 17.1.1 Why Sampling?

There are many reasons that we may wish to draw samples from a probability distribution. Sampling provides a flexible way to approximate many sums and

integrals at reduced cost. Sometimes we use this to provide a significant speedup to a costly but tractable sum, as in the case when we subsample the full training cost with minibatches. In other cases, our learning algorithm requires us to approximate an intractable sum or integral, such as the gradient of the log partition function of an undirected model. In many other cases, sampling is actually our goal, in the sense that we want to train a model that can sample from the training distribution.

### 17.1.2 Basics of Monte Carlo Sampling

When a sum or an integral cannot be computed exactly (for example the sum has an exponential number of terms and no exact simplification is known) it is often possible to approximate it using Monte Carlo sampling. The idea is to view the sum or integral as if it was an expectation under some distribution and to *approximate the expectation by a corresponding average*. Let

$$s = \sum_{\mathbf{x}} p(\mathbf{x}) f(\mathbf{x}) = E_p[f(\mathbf{x})] \quad (17.1)$$

or

$$s = \int p(\mathbf{x}) f(\mathbf{x}) d\mathbf{x} = E_p[f(\mathbf{x})] \quad (17.2)$$

be the sum or integral to estimate, rewritten as an expectation, with the constraint that  $p$  is a probability distribution (for the sum) or a probability density (for the integral) over random variable  $\mathbf{x}$ .

We can approximate  $s$  by drawing  $n$  samples  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$  from  $p$  and then forming the empirical average

$$\hat{s}_n = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}^{(i)}). \quad (17.3)$$

This approximation is justified by a few different properties. The first trivial observation is that the estimator  $\hat{s}$  is unbiased, since

$$\mathbb{E}[\hat{s}_n] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[f(\mathbf{x}^{(i)})] = \frac{1}{n} \sum_{i=1}^n s = s. \quad (17.4)$$

But in addition, the *law of large numbers* states that if the samples  $\mathbf{x}^{(i)}$  are i.i.d., then the average converges almost surely to the expected value:

$$\lim_{n \rightarrow \infty} \hat{s}_n = s, \quad (17.5)$$

provided that the variance of the individual terms,  $\text{Var}[f(\mathbf{x}^{(i)})]$ , is bounded. To see this more clearly, consider the variance of  $\hat{s}_n$  as  $n$  increases. The variance  $\text{Var}[\hat{s}_n]$  decreases and converges to 0, so long as  $\text{Var}[f(\mathbf{x}^{(i)})] < \infty$ :

$$\text{Var}[\hat{s}_n] = \frac{1}{n^2} \sum_{i=1}^n \text{Var}[f(\mathbf{x})] \quad (17.6)$$

$$= \frac{\text{Var}[f(\mathbf{x})]}{n}. \quad (17.7)$$

This convenient result also tells us how to estimate the uncertainty in a Monte Carlo average or equivalently the amount of expected error of the Monte Carlo approximation. We compute both the empirical average of the  $f(\mathbf{x}^{(i)})$  and their empirical variance,<sup>1</sup> and then divide the estimated variance by the number of samples  $n$  to obtain an estimator of  $\text{Var}[\hat{s}_n]$ . The *central limit theorem* tells us that the distribution of the average,  $\hat{s}_n$ , converges to a normal distribution with mean  $s$  and variance  $\frac{\text{Var}[f(\mathbf{x})]}{n}$ . This allows us to estimate confidence intervals around the estimate  $\hat{s}_n$ , using the cumulative distribution of the normal density.

However, all this relies on our ability to easily sample from the base distribution  $p(\mathbf{x})$ , but doing so is not always possible. When it is not feasible to sample from  $p$ , an alternative is to use importance sampling, presented in Sec. 17.2. A more general approach is to form a sequence of estimators that converge towards the distribution of interest. That is the approach of Monte Carlo Markov chains (Sec. 17.3).

## 17.2 Importance Sampling

An important step in the decomposition of the integrand (or summand) used by the Monte Carlo method in Eq. 17.2 is deciding which part of the integrand should play the role the probability  $p(\mathbf{x})$  and which part of the integrand should play the role of the quantity  $f(\mathbf{x})$  whose expected value (under that probability distribution) is to be estimated. There is no unique decomposition because  $p(\mathbf{x})f(\mathbf{x})$  can always be rewritten as

$$p(\mathbf{x})f(\mathbf{x}) = q(\mathbf{x}) \frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})}, \quad (17.8)$$

where we now sample from  $q$  and average  $\frac{pf}{q}$ . In many cases, we wish to compute an expectation for a given  $p$  and an  $f$ , and the fact that the problem is specified

---

<sup>1</sup>The unbiased estimator of the variance is often preferred, in which the sum of squared differences is divided by  $n - 1$  instead of  $n$ .

from the start as an expectation suggests that this  $p$  and  $f$  would be a natural choice of decomposition. However, the original specification of the problem may not be the optimal choice in terms of the number of samples required to obtain a given level of accuracy. Fortunately, the form of the optimal choice  $q^*$  can be derived easily. The optimal  $q^*$  corresponds to what is called optimal importance sampling.

Because of the identity shown in Eq. 17.8, any Monte Carlo estimator

$$\hat{s}_p = \frac{1}{n} \sum_{i=1, \mathbf{x}^{(i)} \sim p}^n f(\mathbf{x}^{(i)}) \quad (17.9)$$

can be transformed into an importance sampling estimator

$$\hat{s}_q = \frac{1}{n} \sum_{i=1, \mathbf{x}^{(i)} \sim q}^n \frac{p(\mathbf{x}^{(i)})f(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})}. \quad (17.10)$$

We see readily that the expected value of the estimator does not depend on  $q$ :

$$\mathbb{E}_q[\hat{s}_q] = \mathbb{E}_q[\hat{s}_p] = s. \quad (17.11)$$

However, the variance of an importance sampling estimator can be greatly sensitive to the choice of  $q$ . The variance is given by

$$\text{Var}[\hat{s}_q] = \text{Var}\left[\frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})}\right]/n. \quad (17.12)$$

The minimum variance occurs when  $q$  is

$$q^*(\mathbf{x}) = \frac{p(\mathbf{x})|f(\mathbf{x})|}{Z}, \quad (17.13)$$

where  $Z$  is the normalization constant, chosen so that  $q^*(\mathbf{x})$  sums or integrates to 1 as appropriate. Better importance sampling distributions put more weight where the integrand is larger. In fact, when  $f(\mathbf{x})$  does not change sign,  $\text{Var}[\hat{s}_{q^*}] = 0$ , meaning that **a single sample is sufficient** when the optimal distribution is used. Of course, this is only because the computation of  $q^*$  has essentially solved the original problem, so it is usually not practical to use this approach of drawing a single sample from the optimal distribution.

Any choice of sampling distribution  $q$  is valid (in the sense of yielding the correct expected value) and  $q^*$  is the optimal one (in the sense of yielding minimum variance). Sampling from  $q^*$  is usually infeasible, but other choices of  $q$  can be feasible while still reducing the variance somewhat.

Another approach is to use **biased importance sampling**, which has the advantage of not requiring normalized  $p$  or  $q$ . In the case of discrete variables, the biased importance sampling estimator is given by

$$\hat{s}_{BIS} = \frac{\sum_{i=1}^n \frac{p(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})} f(\mathbf{x}^{(i)})}{\sum_{i=1}^n \frac{p(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})}} \quad (17.14)$$

$$= \frac{\sum_{i=1}^n \frac{p(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})} f(\mathbf{x}^{(i)})}{\sum_{i=1}^n \frac{p(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})}} \quad (17.15)$$

$$= \frac{\sum_{i=1}^n \frac{\tilde{p}(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})} f(\mathbf{x}^{(i)})}{\sum_{i=1}^n \frac{\tilde{p}(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})}}, \quad (17.16)$$

where  $\tilde{p}$  and  $\tilde{q}$  are the unnormalized forms of  $p$  and  $q$  and the  $\mathbf{x}^{(i)}$  are the samples from  $q$ . This estimator is biased because  $\mathbb{E}[\hat{s}_{BIS}] \neq s$ , except asymptotically when  $n \rightarrow \infty$  and the denominator of Eq. 17.14 converges to 1. Hence this estimator is called asymptotically unbiased.

Although a good choice of  $q$  can greatly improve the efficiency of Monte Carlo estimation, a poor choice of  $q$  can make the efficiency much worse. Going back to Eq. 17.12, we see that if there are samples of  $q$  for which  $\frac{p(\mathbf{x})|f(\mathbf{x})|}{q(\mathbf{x})}$  is large, then the variance of the estimator can get very large. This may happen when  $q(\mathbf{x})$  is tiny while neither  $p(\mathbf{x})$  nor  $f(\mathbf{x})$  are small enough to cancel it. The  $q$  distribution is usually chosen to be a very simple distribution so that it is easy to sample from. When  $\mathbf{x}$  is high-dimensional, this simplicity in  $q$  causes it to match  $p$  or  $p|f|$  poorly. When  $q(\mathbf{x}^{(i)}) \gg p(\mathbf{x}^{(i)})|f(\mathbf{x}^{(i)})|$ , importance sampling collects useless samples (summing tiny numbers or zeros). On the other hand, when  $q(\mathbf{x}^{(i)}) \ll p(\mathbf{x}^{(i)})|f(\mathbf{x}^{(i)})|$ , which will happen more rarely, the ratio can be huge. Because these latter events are rare, they may not show up in a typical sample, yielding typical underestimation of  $s$ , compensated rarely by gross overestimation. Such very large or very small numbers are typical when  $\mathbf{x}$  is high dimensional, because in high dimension the dynamic range of joint probabilities can be very large.

In spite of this danger, importance sampling and its variants have been found very useful in many machine learning algorithms, including deep learning algorithms. For example, see the use of importance sampling to accelerate training in neural language models with a large vocabulary (Sec. 12.4.3.3) or other neural nets with a large number of outputs. See also how importance sampling has been used to estimate a partition function (the normalization constant of a probability

distribution) in Sec. 18.7, and to estimate the log-likelihood in deep directed models such as the variational autoencoder, in Sec. 20.10.3. Importance sampling may also be used to improve the estimate of the gradient of the cost function used to train model parameters with stochastic gradient descent, particularly for models such as classifiers where most of the total value of the cost function comes from a small number of misclassified examples. Sampling more difficult examples more frequently can reduce the variance of the gradient in such cases (Hinton, 2006).

### 17.3 Markov Chain Monte Carlo Methods

In many cases, we wish to use a Monte Carlo technique but there is no tractable method for drawing exact samples from the distribution  $p_{\text{model}}(\mathbf{x})$  or from a good (low variance) importance sampling distribution  $q(\mathbf{x})$ . In the context of deep learning, this most often happens when  $p_{\text{model}}(\mathbf{x})$  is represented by an undirected model. In these cases, we introduce a mathematical tool called a *Markov chain* to approximately sample from  $p_{\text{model}}(\mathbf{x})$ . The family of algorithms that use Markov chains to perform Monte Carlo estimates is called *Markov chain Monte Carlo methods* (MCMC). Markov chain Monte Carlo methods for machine learning are described at greater length in Koller and Friedman (2009). The most standard, generic guarantees for MCMC techniques are only applicable when the model does not assign zero probability to any state. Therefore, it is most convenient to present these techniques as sampling from an energy-based model (EBM)  $p(\mathbf{x}) \propto \exp(-E(\mathbf{x}))$  as described in Sec. 16.2.4. In the EBM formulation, every state is guaranteed to have non-zero probability. MCMC methods are in fact more broadly applicable and can be used with many probability distributions that contain zero probability states. However, the theoretical guarantees concerning the behavior of MCMC methods must be proven on a case-by-case basis for different families of such distributions. In the context of deep learning, it is most common to rely on the most general theoretical guarantees that naturally apply to all energy-based models.

To understand why drawing samples from an energy-based model is difficult, consider an EBM over just two variables, defining a distribution  $p(a, b)$ . In order to sample  $a$ , we must draw  $a$  from  $p(a | b)$ , and in order to sample  $b$ , we must draw  $b$  from  $p(b | a)$ . It seems to be an intractable chicken-and-egg problem. Directed models avoid this because their graph is directed and acyclic. To perform *ancestral sampling* one simply samples each of the variables in topological order, conditioning on each variable’s parents, which are guaranteed to have already been sampled (Sec. 16.3). Ancestral sampling defines an efficient, single-pass method of

obtaining a sample.

In an EBM, we can avoid this chicken and egg problem by sampling using a Markov chain. The core idea of a Markov chain is to have a state  $\mathbf{x}$  that begins as an arbitrary value. Over time, we randomly update  $\mathbf{x}$  repeatedly. Eventually  $\mathbf{x}$  becomes (very nearly) a fair sample from  $p(\mathbf{x})$ . Formally, a Markov chain is defined by a random state  $\mathbf{x}$  and a transition distribution  $T(\mathbf{x}' \mid \mathbf{x})$  specifying the probability that a random update will go to state  $\mathbf{x}'$  if it starts in state  $\mathbf{x}$ . Running the Markov chain means repeatedly updating the state  $\mathbf{x}$  to a value  $\mathbf{x}'$  sampled from  $T(\mathbf{x}' \mid \mathbf{x})$ .

To gain some theoretical understanding of how MCMC methods work, it is useful to reparametrize the problem. First, we restrict our attention to the case where the random variable  $\mathbf{x}$  has countably many states. In this case, we can represent the state as just a positive integer  $x$ . Different integer values of  $x$  map back to different states  $\mathbf{x}$  in the original problem.

Consider what happens when we run infinitely many Markov chains in parallel. All of the states of the different Markov chains are drawn from some distribution  $q^{(t)}(x)$ , where  $t$  indicates the number of time steps that have elapsed. At the beginning,  $q^{(0)}$  is some distribution that we used to arbitrarily initialize  $x$  for each Markov chain. Later,  $q^{(t)}$  is influenced by all of the Markov chain steps that have run so far. Our goal is for  $q^{(t)}(x)$  to converge to  $p(x)$ .

Because we have reparametrized the problem in terms of positive integer  $x$ , we can describe the probability distribution  $q$  using a vector  $\mathbf{v}$ , with

$$q(\mathbf{x} = i) = v_i. \quad (17.17)$$

Consider what happens when we update a single Markov chain's state  $x$  to a new state  $x'$ . The probability of a single state landing in state  $x'$  is given by

$$q^{(t+1)}(x') = \sum_x q^{(t)}(x)T(x' \mid x). \quad (17.18)$$

Using our integer parametrization, we can represent the effect of the transition operator  $T$  using a matrix  $\mathbf{A}$ . We define  $\mathbf{A}$  so that

$$A_{i,j} = T(\mathbf{x}' = i \mid \mathbf{x} = j). \quad (17.19)$$

Using this definition, we can now rewrite Eq. 17.18. Rather than writing it in terms of  $q$  and  $T$  to understand how a single state is updated, we may now use  $\mathbf{v}$  and  $\mathbf{A}$  to describe how the entire distribution over all the different Markov chains run in parallel shifts as we apply an update:

$$\mathbf{v}^{(t)} = \mathbf{A}\mathbf{v}^{(t-1)}. \quad (17.20)$$

Applying the Markov chain update repeatedly corresponds to multiplying by the matrix  $\mathbf{A}$  repeatedly. In other words, we can think of the process as exponentiating the matrix  $\mathbf{A}$ :

$$\mathbf{v}^{(t)} = \mathbf{A}^t \mathbf{v}^{(0)}. \quad (17.21)$$

The matrix  $\mathbf{A}$  has special structure because each of its columns represents a probability distribution. Such matrices are called *stochastic matrices*. If there is a non-zero probability of transitioning from any state  $x$  to any other state  $x'$  for some power  $t$ , then the Perron-Frobenius theorem (Perron, 1907; Frobenius, 1908) guarantees that the largest eigenvalue is real and equal to 1. Over time, we can see that all of the eigenvalues are exponentiated:

$$\mathbf{v}^{(t)} = (\mathbf{V} \text{diag}(\boldsymbol{\lambda}) \mathbf{V}^{-1})^t \mathbf{v}^{(0)} = \mathbf{V} \text{diag}(\boldsymbol{\lambda})^t \mathbf{V}^{-1} \mathbf{v}^{(0)}. \quad (17.22)$$

This process causes all of the eigenvalues that are not equal to 1 to decay to zero. Under some additional mild conditions,  $\mathbf{A}$  is guaranteed to have only one eigenvector with eigenvalue 1. The process thus converges to a *stationary distribution*, sometimes also called the *equilibrium distribution*. At convergence,

$$\mathbf{v}' = \mathbf{A}\mathbf{v} = \mathbf{v}, \quad (17.23)$$

and this same condition holds for every additional step. This is an eigenvector equation. To be a stationary point,  $\mathbf{v}$  must be an eigenvector with corresponding eigenvalue 1. This condition guarantees that once we have reached the stationary distribution, repeated applications of the transition sampling procedure do not change the *distribution* over the states of all the various Markov chains (although transition operator does change each individual state, of course).

If we have chosen  $T$  correctly, then the stationary distribution  $q$  will be equal to the distribution  $p$  we wish to sample from. We will describe how to choose  $T$  shortly, in Sec. 17.4.

Most properties of Markov Chains with countable states can be generalized to continuous variables. In this situation, some authors call the Markov Chain a *Harris chain* but we use the term Markov Chain to describe both conditions. In general, a Markov chain with transition operator  $T$  will converge, under mild conditions, to a fixed point described by the equation

$$q'(\mathbf{x}') = \mathbb{E}_{\mathbf{x} \sim q} T(\mathbf{x}' | \mathbf{x}), \quad (17.24)$$

which in the discrete case is just rewriting Eq. 17.23. When  $\mathbf{x}$  is discrete, the expectation corresponds to a sum, and when  $\mathbf{x}$  is continuous, the expectation corresponds to an integral.

Regardless of whether the state is continuous or discrete, all Markov chain methods consist of repeatedly applying stochastic updates until eventually the state begins to yield samples from the equilibrium distribution. Running the Markov chain until it reaches its equilibrium distribution is called “*burning in*” the Markov chain. After the chain has reached equilibrium, a sequence of infinitely many samples may be drawn from the equilibrium distribution. They are identically distributed but any two successive samples will be highly correlated with each other. A finite sequence of samples may thus not be very representative of the equilibrium distribution. One way to mitigate this problem is to return only every  $n$  successive samples, so that our estimate of the statistics of the equilibrium distribution is not as biased by the correlation between an MCMC sample and the next several samples. Markov chains are thus expensive to use because of the time required to burn in to the equilibrium distribution and the time required to transition from one sample to another reasonably decorrelated sample after reaching equilibrium. If one desires truly independent samples, one can run multiple Markov chains in parallel. This approach uses extra parallel computation to eliminate latency. The strategy of using only a single Markov chain to generate all samples and the strategy of using one Markov chain for each desired sample are two extremes; deep learning practitioners usually use a number of chains that is similar to the number of examples in a minibatch and then draw as many samples as are needed from this fixed set of Markov chains. A commonly used number of Markov chains is 100.

Another difficulty is that we do not know in advance how many steps the Markov chain must run before reaching its equilibrium distribution. This length of time is called the *mixing time*. It is also very difficult to test whether a Markov chain has reached equilibrium. We do not have a precise enough theory for guiding us in answering this question. Theory tells us that the chain will converge, but not much more. If we analyze the Markov chain from the point of view of a matrix  $\mathbf{A}$  acting on a vector of probabilities  $\mathbf{v}$ , then we know that the chain mixes when  $\mathbf{A}^t$  has effectively lost all of the eigenvalues from  $\mathbf{A}$  besides the unique eigenvalue of 1. This means that the magnitude of the second largest eigenvalue will determine the mixing time. However, in practice, we cannot actually represent our Markov chain in terms of a matrix. The number of states that our probabilistic model can visit is exponentially large in the number of variables, so it is infeasible to represent  $\mathbf{v}$ ,  $\mathbf{A}$ , or the eigenvalues of  $\mathbf{A}$ . Due to these and other obstacles, we usually do not know whether a Markov chain has mixed. Instead, we simply run the Markov chain for an amount of time that we roughly estimate to be sufficient, and use heuristic methods to determine whether the chain has mixed. These heuristic methods include manually inspecting samples or measuring correlations between successive samples.

## 17.4 Gibbs Sampling

So far we have described how to draw samples from a distribution  $q(\mathbf{x})$  by repeatedly updating  $\mathbf{x} \leftarrow \mathbf{x}' \sim T(\mathbf{x}' | \mathbf{x})$ . However, we have not described how to ensure that  $q(\mathbf{x})$  is a useful distribution. Two basic approaches are considered in this book. The first one is to derive  $T$  from a given learned  $p_{\text{model}}$ , described below with the case of sampling from EBMs. The second one is to directly parametrize  $T$  and learn it, so that its stationary distribution implicitly defines the  $p_{\text{model}}$  of interest. Examples of this second approach are discussed in Sec. 20.12 and Sec. 20.13.

In the context of deep learning, we commonly use Markov chains to draw samples from an energy-based model defining a distribution  $p_{\text{model}}(\mathbf{x})$ . In this case, we want the  $q(\mathbf{x})$  for the Markov chain to be  $p_{\text{model}}(\mathbf{x})$ . To obtain the desired  $q(\mathbf{x})$ , we must choose an appropriate  $T(\mathbf{x}' | \mathbf{x})$ .

A conceptually simple and effective approach to building a Markov chain that samples from  $p_{\text{model}}(\mathbf{x})$  is to use *Gibbs sampling*, in which sampling from  $T(\mathbf{x}' | \mathbf{x})$  is accomplished by selecting one variable  $x_i$  and sampling it from  $p_{\text{model}}$  conditioned on its neighbors in the undirected graph  $\mathcal{G}$  defining the structure of the energy-based model. It is also possible to sample several variables at the same time so long as they are conditionally independent given all of their neighbors. As shown in the RBM example in Sec. 16.7.1, all of the hidden units of an RBM may be sampled simultaneously because they are conditionally independent from each other given all of the visible units. Likewise, all of the visible units may be sampled simultaneously because they are conditionally independent from each other given all of the hidden units. Gibbs sampling approaches that update many variables simultaneously in this way are called *block Gibbs sampling*.

Alternate approaches to designing Markov chains to sample from  $p_{\text{model}}$  are possible. For example, the Metropolis-Hastings algorithm is widely used in other disciplines. In the context of the deep learning approach to undirected modeling, it is rare to use any approach other than Gibbs sampling. Improved sampling techniques are one possible research frontier.

## 17.5 The Challenge of Mixing between Separated Modes

The primary difficulty involved with MCMC methods is that they have a tendency to *mix* poorly. Ideally, successive samples from a Markov chain designed to sample from  $p(\mathbf{x})$  would be completely independent from each other and would visit many different regions in  $\mathbf{x}$  space proportional to their probability. Instead, especially in high dimensional cases, MCMC samples become very correlated. We refer

to such behavior as slow mixing or even failure to mix. MCMC methods with slow mixing can be seen as inadvertently performing something resembling noisy gradient descent on the energy function, or equivalently noisy hill climbing on the probability, with respect to the state of the chain (the random variables being sampled). The chain tends to take small steps (in the space of the state of the Markov chain), from a configuration  $\mathbf{x}^{(t-1)}$  to a configuration  $\mathbf{x}^{(t)}$ , with the energy  $E(\mathbf{x}^{(t)})$  generally lower or approximately equal to the energy  $E(\mathbf{x}^{(t-1)})$ , with a preference for moves that yield lower energy configurations. When starting from a rather improbable configuration (higher energy than the typical ones from  $p(\mathbf{x})$ ), the chain tends to gradually reduce the energy of the state and only occasionally move to another mode. Once the chain has found a region of low energy (for example, if the variables are pixels in an image, a region of low energy might be a connected manifold of images of the same object), which we call a mode, the chain will tend to walk around that mode (following a kind of random walk). Once in a while it will step out of that mode and generally return to it or (if it finds an escape route) move towards another mode. The problem is that successful escape routes are rare for many interesting distributions, so the Markov chain will continue to sample the same mode longer than it should.

This is very clear when we consider the Gibbs sampling algorithm (Sec. 17.4). In this context, consider the probability of going from one mode to a nearby mode within a given number of steps. What will determine that probability is the shape of the “energy barrier” between these modes. Transitions between two modes that are separated by a high energy barrier (a region of low probability) are exponentially less likely (in terms of the height of the energy barrier). This is illustrated in Fig. 17.1. The problem arises when there are multiple modes with high probability that are separated by regions of low probability, especially when each Gibbs sampling step must update only a small subset of variables whose values are largely determined by the other variables.

As a simple example, consider an energy-based model over two variables  $a$  and  $b$ , which are both binary with a sign, taking on values  $-1$  and  $1$ . If  $E(a, b) = -wab$  for some large positive number  $w$ , then the model expresses a strong belief that  $a$  and  $b$  have the same sign. Consider updating  $b$  using a Gibbs sampling step with  $a = 1$ . The conditional distribution over  $b$  is given by  $P(b = 1 | a = 1) = \sigma(w)$ . If  $w$  is large, the sigmoid saturates, and the probability of also assigning  $b$  to be  $1$  is close to  $1$ . Likewise, if  $a = -1$ , the probability of assigning  $b$  to be  $-1$  is close to  $1$ . According to  $P_{\text{model}}(a, b)$ , both signs of both variables are equally likely. According to  $P_{\text{model}}(a | b)$ , both variables should have the same sign. This means that Gibbs sampling will only very rarely flip the signs of these variables.

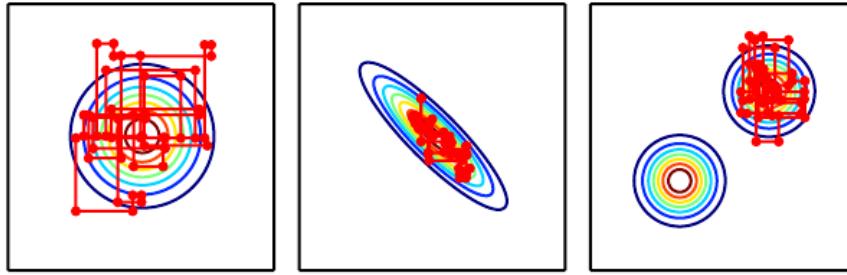


Figure 17.1: Paths followed by Gibbs sampling for three distributions, with the Markov chain initialized at the mode in both cases. (*Left*) A multivariate normal distribution with two independent variables. Gibbs sampling mixes well because the variables are independent. (*Center*) A multivariate normal distribution with highly correlated variables. The correlation between variables makes it difficult for the Markov chain to mix. Because each variable must be updated conditioned on the other, the correlation reduces the rate at which the Markov chain can move away from the starting point. (*Right*) A mixture of Gaussians with widely separated modes that are not axis-aligned. Gibbs sampling mixes very slowly because it is difficult to change modes while altering only one variable at a time.

In more practical scenarios, the challenge is even greater because we care not only about making transitions between two modes but more generally between all the many modes that a real model might contain. If several such transitions are difficult because of the difficulty of mixing between modes, then it becomes very expensive to obtain a reliable set of samples covering most of the modes, and convergence of the chain to its stationary distribution is very slow.

Sometimes this problem can be resolved by finding groups of highly dependent units and updating all of them simultaneously in a block. Unfortunately, when the dependencies are complicated, it can be computationally intractable to draw a sample from the group. After all, the problem that the Markov chain was originally introduced to solve is this problem of sampling from a large group of variables.

In the context of models with latent variables, which define a joint distribution  $p_{\text{model}}(\mathbf{x}, \mathbf{h})$ , we often draw samples of  $\mathbf{x}$  by alternating between sampling from  $p_{\text{model}}(\mathbf{x} | \mathbf{h})$  and sampling from  $p_{\text{model}}(\mathbf{h} | \mathbf{x})$ . From the point of view of mixing rapidly, we would like  $p_{\text{model}}(\mathbf{h} | \mathbf{x})$  to have very high entropy. However, from the point of view of learning a useful representation of  $\mathbf{h}$ , we would like  $\mathbf{h}$  to encode

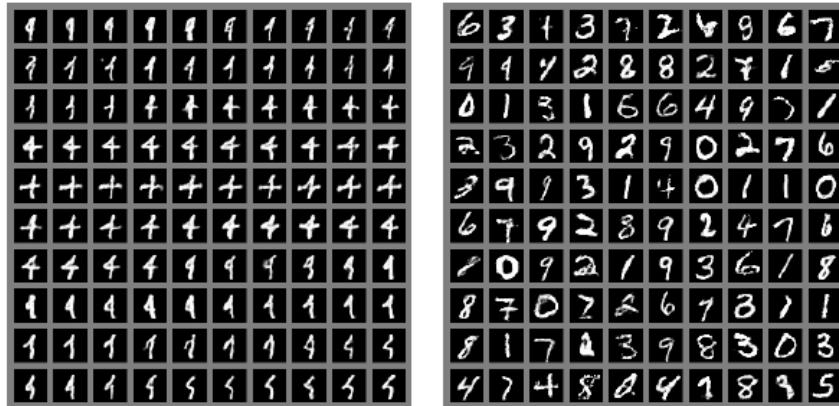


Figure 17.2: An illustration of the slow mixing problem in deep probabilistic models. Each panel should be read left to right, top to bottom. (*Left*) Consecutive samples from Gibbs sampling applied to a deep Boltzmann machine trained on the MNIST dataset. Consecutive samples are similar to each other. Because the Gibbs sampling is performed in a deep graphical model, this similarity is based more on semantic rather than raw visual features, but it is still difficult for the Gibbs chain to transition from one mode of the distribution to another, for example by changing the digit identity. (*Right*) Consecutive ancestral samples from a generative adversarial network. Because ancestral sampling generates each sample independently from the others, there is no mixing problem.

enough information about  $\mathbf{x}$  to reconstruct it well, which implies that  $\mathbf{h}$  and  $\mathbf{x}$  should have very high mutual information. These two goals are at odds with each other. We often learn generative models that very precisely encode  $\mathbf{x}$  into  $\mathbf{h}$  but are not able to mix very well. This situation arises frequently with Boltzmann machines—the sharper the distribution a Boltzmann machine learns, the harder it is for a Markov chain sampling from the model distribution to mix well. This problem is illustrated in Fig. 17.2.

All this could make MCMC methods less useful when the distribution of interest has a manifold structure with a separate manifold for each class: the distribution is concentrated around many modes and these modes are separated by vast regions of high energy. This type of distribution is what we expect in many classification problems and would make MCMC methods converge very slowly because of poor mixing between modes.

### 17.5.1 Tempering to Mix between Modes

When a distribution has sharp peaks of high probability surrounded by regions of low probability, it is difficult to mix between the different modes of the distribution.

Several techniques for faster mixing are based on constructing alternative versions of the target distribution in which the peaks are not as high and the surrounding valleys are not as low. Energy-based models provide a particularly simple way to do so. So far, we have described an energy-based model as defining a probability distribution

$$p(\mathbf{x}) \propto \exp(-E(\mathbf{x})). \quad (17.25)$$

Energy-based models may be augmented with an extra parameter  $\beta$  controlling how sharply peaked the distribution is:

$$p_\beta(\mathbf{x}) \propto \exp(-\beta E(\mathbf{x})). \quad (17.26)$$

The  $\beta$  parameter is often described as being the reciprocal of the *temperature*, reflecting the origin of energy-based models in statistical physics. When the temperature falls to zero and  $\beta$  rises to infinity, the energy-based model becomes deterministic. When the temperature rises to infinity and  $\beta$  falls to zero, the distribution (for discrete  $\mathbf{x}$ ) becomes uniform.

Typically, a model is trained to be evaluated at  $\beta = 1$ . However, we can make use of other temperatures, particularly those where  $\beta < 1$ . *Tempering* is a general strategy of mixing between modes of  $p_1$  rapidly by drawing samples with  $\beta < 1$ .

Markov chains based on *tempered transitions* (Neal, 1994) temporarily sample from higher-temperature distributions in order to mix to different modes, then resume sampling from the unit temperature distribution. These techniques have been applied to models such as RBMs (Salakhutdinov, 2010). Another approach is to use *parallel tempering* (Iba, 2001), in which the Markov chain simulates many different states in parallel, at different temperatures. The highest temperature states mix slowly, while the lowest temperature states, at temperature 1, provide accurate samples from the model. The transition operator includes stochastically swapping states between two different temperature levels, so that a sufficiently high-probability sample from a high-temperature slot can jump into a lower temperature slot. This approach has also been applied to RBMs (Desjardins *et al.*, 2010; Cho *et al.*, 2010). Although tempering is a promising approach, at this point it has not allowed researchers to make a strong advance in solving the challenge of sampling from complex EBMs. One possible reason is that there are *critical temperatures* around which the temperature transition must be very slow (as the temperature is gradually reduced) in order for tempering to be effective.

### 17.5.2 Depth May Help Mixing

When drawing samples from a latent variable model  $p(\mathbf{h}, \mathbf{x})$ , we have seen that if  $p(\mathbf{h} | \mathbf{x})$  encodes  $\mathbf{x}$  too well, then sampling from  $p(\mathbf{x} | \mathbf{h})$  will not change  $\mathbf{x}$  very

much and mixing will be poor. One way to resolve this problem is to make  $\mathbf{h}$  be a deep representation, that encodes  $\mathbf{x}$  into  $\mathbf{h}$  in such a way that a Markov chain in the space of  $\mathbf{h}$  can mix more easily. Many representation learning algorithms, such as autoencoders and RBMs, tend to yield a marginal distribution over  $\mathbf{h}$  that is more uniform and more unimodal than the original data distribution over  $\mathbf{x}$ . It can be argued that this arises from trying to minimize reconstruction error while using all of the available representation space, because minimizing reconstruction error over the training examples will be better achieved when different training examples are easily distinguishable from each other in  $\mathbf{h}$ -space, and thus well separated. [Bengio et al. \(2013a\)](#) observed that deeper stacks of regularized autoencoders or RBMs yield marginal distributions in the top-level  $\mathbf{h}$ -space that appeared more spread out and more uniform, with less of a gap between the regions corresponding to different modes (categories, in the experiments). Training an RBM in that higher-level space allowed Gibbs sampling to mix faster between modes. It remains however unclear how to exploit this observation to help better train and sample from deep generative models.

Despite the difficulty of mixing, Monte Carlo techniques are useful and are often the best tool available. Indeed, they are the primary tool used to confront the intractable partition function of undirected models, discussed next.

## Chapter 18

# Confronting the Partition Function

In Sec. 16.2.2 we saw that many probabilistic models (commonly known as undirected graphical models) are defined by an unnormalized probability distribution  $\tilde{p}(\mathbf{x}; \boldsymbol{\theta})$ . We must normalize  $\tilde{p}$  by dividing by a partition function  $Z(\boldsymbol{\theta})$  in order to obtain a valid probability distribution:

$$p(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \tilde{p}(\mathbf{x}; \boldsymbol{\theta}). \quad (18.1)$$

The partition function is an integral (for continuous variables) or sum (for discrete variables) over the unnormalized probability of all states:

$$\int \tilde{p}(\mathbf{x}) d\mathbf{x} \quad (18.2)$$

or

$$\sum_{\mathbf{x}} \tilde{p}(\mathbf{x}). \quad (18.3)$$

This operation is intractable for many interesting models.

As we will see in Chapter 20, several deep learning models are designed to have a tractable normalizing constant, or are designed to be used in ways that do not involve computing  $p(\mathbf{x})$  at all. However, other models directly confront the challenge of intractable partition functions. In this chapter, we describe techniques used for training and evaluating models that have intractable partition functions.

## 18.1 The Log-Likelihood Gradient

What makes learning undirected models by maximum likelihood particularly difficult is that the partition function depends on the parameters. The gradient of the log-likelihood with respect to the parameters has a term corresponding to the gradient of the partition function:

$$\nabla_{\boldsymbol{\theta}} \log p(\mathbf{x}; \boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}; \boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}} \log Z(\boldsymbol{\theta}). \quad (18.4)$$

This is a well-known decomposition into the *positive phase* and *negative phase* of learning.

For most undirected models of interest, the negative phase is difficult. Models with no latent variables or with few interactions between latent variables typically have a tractable positive phase. The quintessential example of a model with a straightforward positive phase and difficult negative phase is the RBM, which has hidden units that are conditionally independent from each other given the visible units. The case where the positive phase is difficult, with complicated interactions between latent variables, is primarily covered in Chapter 19. This chapter focuses on the difficulties of the negative phase.

Let us look more closely at the gradient of  $\log Z$ :

$$\frac{\partial}{\partial \boldsymbol{\theta}} \log Z \quad (18.5)$$

$$= \frac{\frac{\partial}{\partial \boldsymbol{\theta}} Z}{Z} \quad (18.6)$$

$$= \frac{\frac{\partial}{\partial \boldsymbol{\theta}} \sum_{\mathbf{x}} \tilde{p}(\mathbf{x})}{Z} \quad (18.7)$$

$$= \frac{\sum_{\mathbf{x}} \frac{\partial}{\partial \boldsymbol{\theta}} \tilde{p}(\mathbf{x})}{Z}. \quad (18.8)$$

For models that guarantee  $p(\mathbf{x}) > 0$  for all  $\mathbf{x}$ , we can substitute  $\exp(\log \tilde{p}(\mathbf{x}))$  for  $\tilde{p}(\mathbf{x})$ :

$$\frac{\sum_{\mathbf{x}} \frac{\partial}{\partial \boldsymbol{\theta}} \exp(\log \tilde{p}(\mathbf{x}))}{Z} \quad (18.9)$$

$$= \frac{\sum_{\mathbf{x}} \exp(\log \tilde{p}(\mathbf{x})) \frac{\partial}{\partial \boldsymbol{\theta}} \log \tilde{p}(\mathbf{x})}{Z} \quad (18.10)$$

$$= \frac{\sum_{\mathbf{x}} \tilde{p}(\mathbf{x}) \frac{\partial}{\partial \boldsymbol{\theta}} \log \tilde{p}(\mathbf{x})}{Z} \quad (18.11)$$

$$= \sum_{\mathbf{x}} p(\mathbf{x}) \frac{\partial}{\partial \boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}) \quad (18.12)$$

$$= \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \frac{\partial}{\partial \boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}). \quad (18.13)$$

This derivation made use of summation over discrete  $\mathbf{x}$ , but a similar result applies using integration over continuous  $\mathbf{x}$ . In the continuous version of the derivation, we use Leibniz's rule for differentiation under the integral sign to obtain the identity

$$\frac{\partial}{\partial \boldsymbol{\theta}} \int \tilde{p}(\mathbf{x}) d\mathbf{x} = \int \frac{\partial}{\partial \boldsymbol{\theta}} \tilde{p}(\mathbf{x}) d\mathbf{x}. \quad (18.14)$$

This identity is only applicable under certain regularity conditions on  $\tilde{p}$  and  $\frac{\partial}{\partial \boldsymbol{\theta}} \tilde{p}(\mathbf{x})$ . In measure theoretic terms, the conditions are: (i)  $\tilde{p}$  must be a Lebesgue-integrable function of  $\mathbf{x}$  for every value of  $\boldsymbol{\theta}$ ; (ii)  $\frac{\partial}{\partial \boldsymbol{\theta}} \tilde{p}(\mathbf{x})$  must exist for all  $\boldsymbol{\theta}$  and almost all  $\mathbf{x}$ ; (iii) There must exist an integrable function  $R(\mathbf{x})$  that bounds  $\frac{\partial}{\partial \boldsymbol{\theta}} \tilde{p}(\mathbf{x})$  (i.e.  $|\frac{\partial}{\partial \boldsymbol{\theta}} \tilde{p}(\mathbf{x})| \leq R(\mathbf{x})$  for all  $\boldsymbol{\theta}$  and almost all  $\mathbf{x}$ ). Fortunately, most machine learning models of interest have these properties.

This identity

$$\nabla_{\boldsymbol{\theta}} \log Z = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}) \quad (18.15)$$

is the basis for a variety of Monte Carlo methods for approximately maximizing the likelihood of models with intractable partition functions.

The Monte Carlo approach to learning undirected models provides an intuitive framework in which we can think of both the positive phase and the negative phase. In the positive phase, we increase  $\log \tilde{p}(\mathbf{x})$  for  $\mathbf{x}$  drawn from the data. In the negative phase, we decrease the partition function by decreasing  $\log \tilde{p}(\mathbf{x})$  drawn from the model distribution.

In the deep learning literature, it is common to parametrize  $\log \tilde{p}$  in terms of an energy function (Eq. 16.7). In this case, we can interpret the positive phase as pushing down on the energy of training examples and the negative phase as pushing up on the energy of samples drawn from the model, as illustrated in Fig. 18.1.

## 18.2 Stochastic Maximum Likelihood and Contrastive Divergence

The naive way of implementing Eq. 18.15 is to compute it by burning in a set of Markov chains from a random initialization every time the gradient is needed.

When learning is performed using stochastic gradient descent, this means the chains must be burned in once per gradient step. This approach leads to the training procedure presented in Algorithm 18.1. The high cost of burning in the Markov chains in the inner loop makes this procedure computationally infeasible, but this procedure is the starting point that other more practical algorithms aim to approximate.

---

**Algorithm 18.1** A naive MCMC algorithm for maximizing the log-likelihood with an intractable partition function using gradient ascent.

---

Set  $\epsilon$ , the step size, to a small positive number.  
Set  $k$ , the number of Gibbs steps, high enough to allow burn in. Perhaps 100 to train an RBM on a small image patch.  
**while** not converged **do**  
    Sample a minibatch of  $m$  examples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  from the training set.  
     $\mathbf{g} \leftarrow \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}^{(i)}; \boldsymbol{\theta})$ .  
    Initialize a set of  $m$  samples  $\{\tilde{\mathbf{x}}^{(1)}, \dots, \tilde{\mathbf{x}}^{(m)}\}$  to random values (e.g., from a uniform or normal distribution, or possibly a distribution with marginals matched to the model's marginals).  
    **for**  $i = 1$  to  $k$  **do**  
        **for**  $j = 1$  to  $m$  **do**  
             $\tilde{\mathbf{x}}^{(j)} \leftarrow \text{gibbs\_update}(\tilde{\mathbf{x}}^{(j)})$ .  
        **end for**  
    **end for**  
     $\mathbf{g} \leftarrow \mathbf{g} - \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\tilde{\mathbf{x}}^{(i)}; \boldsymbol{\theta})$ .  
     $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \epsilon \mathbf{g}$ .  
**end while**

---

We can view the MCMC approach to maximum likelihood as trying to achieve balance between two forces, one pushing up on the model distribution where the data occurs, and another pushing down on the model distribution where the model samples occur. Fig. 18.1 illustrates this process. The two forces correspond to maximizing  $\log \tilde{p}$  and minimizing  $\log Z$ . Several approximations to the negative phase are possible. Each of these approximations can be understood as making the negative phase computationally cheaper but also making it push down in the wrong locations.

Because the negative phase involves drawing samples from the model's distribution, we can think of it as finding points that the model believes in strongly. Because the negative phase acts to reduce the probability of those points, they are generally considered to represent the model's incorrect beliefs about the world.

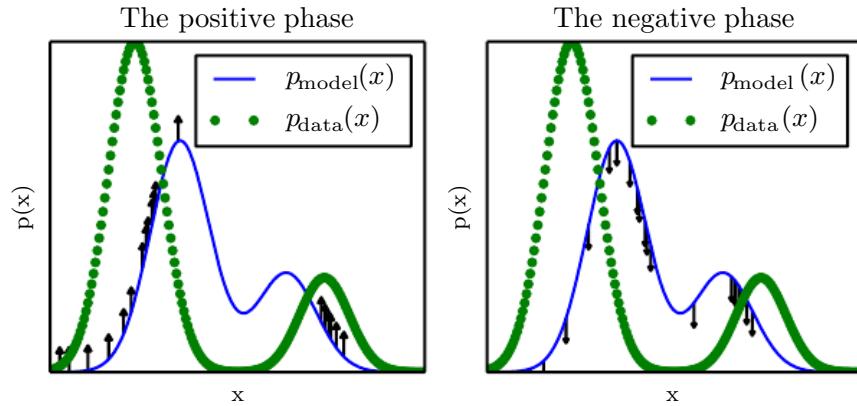


Figure 18.1: The view of Algorithm 18.1 as having a “positive phase” and “negative phase.” (*Left*) In the positive phase, we sample points from the data distribution, and push up on their unnormalized probability. This means points that are likely in the data get pushed up on more. (*Right*) In the negative phase, we sample points from the model distribution, and push down on their unnormalized probability. This counteracts the positive phase’s tendency to just add a large constant to the unnormalized probability everywhere. When the data distribution and the model distribution are equal, the positive phase has the same chance to push up at a point as the negative phase has to push down. When this occurs, there is no longer any gradient (in expectation) and training must terminate.

They are frequently referred to in the literature as “hallucinations” or “fantasy particles.” In fact, the negative phase has been proposed as a possible explanation for dreaming in humans and other animals (Crick and Mitchison, 1983), the idea being that the brain maintains a probabilistic model of the world and follows the gradient of  $\log \tilde{p}$  while experiencing real events while awake and follows the negative gradient of  $\log \tilde{p}$  to minimize  $\log Z$  while sleeping and experiencing events sampled from the current model. This view explains much of the language used to describe algorithms with a positive and negative phase, but it has not been proven to be correct with neuroscientific experiments. In machine learning models, it is usually necessary to use the positive and negative phase simultaneously, rather than in separate time periods of wakefulness and REM sleep. As we will see in Sec. 19.5, other machine learning algorithms draw samples from the model distribution for other purposes and such algorithms could also provide an account for the function of dream sleep.

Given this understanding of the role of the positive and negative phase of learning, we can attempt to design a less expensive alternative to Algorithm 18.1. The main cost of the naive MCMC algorithm is the cost of burning in the Markov

chains from a random initialization at each step. A natural solution is to initialize the Markov chains from a distribution that is very close to the model distribution, so that the burn in operation does not take as many steps.

The *contrastive divergence* (CD, or CD- $k$  to indicate CD with  $k$  Gibbs steps) algorithm initializes the Markov chain at each step with samples from the data distribution (Hinton, 2000, 2010). This approach is presented as Algorithm 18.2. Obtaining samples from the data distribution is free, because they are already available in the data set. Initially, the data distribution is not close to the model distribution, so the negative phase is not very accurate. Fortunately, the positive phase can still accurately increase the model's probability of the data. After the positive phase has had some time to act, the model distribution is closer to the data distribution, and the negative phase starts to become accurate.

---

**Algorithm 18.2** The contrastive divergence algorithm, using gradient ascent as the optimization procedure.

---

```
Set  $\epsilon$ , the step size, to a small positive number.  
Set  $k$ , the number of Gibbs steps, high enough to allow a Markov chain sampling  
from  $p(\mathbf{x}; \boldsymbol{\theta})$  to mix when initialized from  $p_{\text{data}}$ . Perhaps 1-20 to train an RBM  
on a small image patch.  
while not converged do  
    Sample a minibatch of  $m$  examples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  from the training set.  
     $\mathbf{g} \leftarrow \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}^{(i)}; \boldsymbol{\theta}).$   
    for  $i = 1$  to  $m$  do  
         $\tilde{\mathbf{x}}^{(i)} \leftarrow \mathbf{x}^{(i)}.$   
    end for  
    for  $i = 1$  to  $k$  do  
        for  $j = 1$  to  $m$  do  
             $\tilde{\mathbf{x}}^{(j)} \leftarrow \text{gibbs\_update}(\tilde{\mathbf{x}}^{(j)}).$   
        end for  
    end for  
     $\mathbf{g} \leftarrow \mathbf{g} - \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\tilde{\mathbf{x}}^{(i)}; \boldsymbol{\theta}).$   
     $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \epsilon \mathbf{g}.$   
end while
```

---

Of course, CD is still an approximation to the correct negative phase. The main way that CD qualitatively fails to implement the correct negative phase is that it fails to suppress regions of high probability that are far from actual training examples. These regions that have high probability under the model but low probability under the data generating distribution are called *spurious modes*.

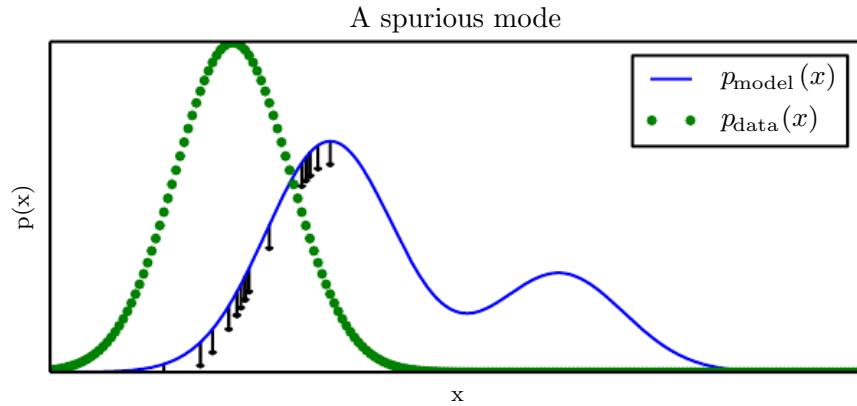


Figure 18.2: An illustration of how the negative phase of contrastive divergence (Algorithm 18.2) can fail to suppress spurious modes. A spurious mode is a mode that is present in the model distribution but absent in the data distribution. Because contrastive divergence initializes its Markov chains from data points and runs the Markov chain for only a few steps, it is unlikely to visit modes in the model that are far from the data points. This means that when sampling from the model, we will sometimes get samples that do not resemble the data. It also means that due to wasting some of its probability mass on these modes, the model will struggle to place high probability mass on the correct modes. For the purpose of visualization, this figure uses a somewhat simplified concept of distance—the spurious mode is far from the correct mode along the number line in  $\mathbb{R}$ . This corresponds to a Markov chain based on making local moves with a single  $x$  variable in  $\mathbb{R}$ . For most deep probabilistic models, the Markov chains are based on Gibbs sampling and can make non-local moves of individual variables but cannot move all of the variables simultaneously. For these problems, it is usually better to consider the edit distance between modes, rather than the Euclidean distance. However, edit distance in a high dimensional space is difficult to depict in a 2-D plot.

Fig. 18.2 illustrates why this happens. Essentially, it is because modes in the model distribution that are far from the data distribution will not be visited by Markov chains initialized at training points, unless  $k$  is very large.

[Carreira-Perpiñan and Hinton \(2005\)](#) showed experimentally that the CD estimator is biased for RBMs and fully visible Boltzmann machines, in that it converges to different points than the maximum likelihood estimator. They argue that because the bias is small, CD could be used as an inexpensive way to initialize a model that could later be fine-tuned via more expensive MCMC methods. [Bengio and Delalleau \(2009\)](#) showed that CD can be interpreted as discarding the smallest terms of the correct MCMC update gradient, which explains the bias.

CD is useful for training shallow models like RBMs. These can in turn be

stacked to initialize deeper models like DBNs or DBMs. However, CD does not provide much help for training deeper models directly. This is because it is difficult to obtain samples of the hidden units given samples of the visible units. Since the hidden units are not included in the data, initializing from training points cannot solve the problem. Even if we initialize the visible units from the data, we will still need to burn in a Markov chain sampling from the distribution over the hidden units conditioned on those visible samples.

The CD algorithm can be thought of as penalizing the model for having a Markov chain that changes the input rapidly when the input comes from the data. This means training with CD somewhat resembles autoencoder training. Even though CD is more biased than some of the other training methods, it can be useful for pretraining shallow models that will later be stacked. This is because the earliest models in the stack are encouraged to copy more information up to their latent variables, thereby making it available to the later models. This should be thought of more of as an often-exploitable side effect of CD training rather than a principled design advantage.

Sutskever and Tieleman (2010) showed that the CD update direction is not the gradient of any function. This allows for situations where CD could cycle forever, but in practice this is not a serious problem.

A different strategy that resolves many of the problems with CD is to initialize the Markov chains at each gradient step with their states from the previous gradient step. This approach was first discovered under the name *stochastic maximum likelihood* (SML) in the applied mathematics and statistics community (Younes, 1998) and later independently rediscovered under the name *persistent contrastive divergence* (PCD, or PCD- $k$  to indicate the use of  $k$  Gibbs steps per update) in the deep learning community (Tieleman, 2008). See Algorithm 18.3. The basic idea of this approach is that, so long as the steps taken by the stochastic gradient algorithm are small, then the model from the previous step will be similar to the model from the current step. It follows that the samples from the previous model’s distribution will be very close to being fair samples from the current model’s distribution, so a Markov chain initialized with these samples will not require much time to mix.

Because each Markov chain is continually updated throughout the learning process, rather than restarted at each gradient step, the chains are free to wander far enough to find all of the model’s modes. SML is thus considerably more resistant to forming models with spurious modes than CD is. Moreover, because it is possible to store the state of all of the sampled variables, whether visible or latent, SML provides an initialization point for both the hidden and visible units.

CD is only able to provide an initialization for the visible units, and therefore requires burn-in for deep models. SML is able to train deep models efficiently. [Marlin et al. \(2010\)](#) compared SML to many of the other criteria presented in this chapter. They found that SML results in the best test set log-likelihood for an RBM, and that if the RBM's hidden units are used as features for an SVM classifier, SML results in the best classification accuracy.

SML is vulnerable to becoming inaccurate if the stochastic gradient algorithm can move the model faster than the Markov chain can mix between steps. This can happen if  $k$  is too small or  $\epsilon$  is too large. The permissible range of values is unfortunately highly problem-dependent. There is no known way to test formally whether the chain is successfully mixing between steps. Subjectively, if the learning rate is too high for the number of Gibbs steps, the human operator will be able to observe that there is much more variance in the negative phase samples across gradient steps rather than across different Markov chains. For example, a model trained on MNIST might sample exclusively 7s on one step. The learning process will then push down strongly on the mode corresponding to 7s, and the model might sample exclusively 9s on the next step.

---

**Algorithm 18.3** The stochastic maximum likelihood / persistent contrastive divergence algorithm using gradient ascent as the optimization procedure.

---

Set  $\epsilon$ , the step size, to a small positive number.

Set  $k$ , the number of Gibbs steps, high enough to allow a Markov chain sampling from  $p(\mathbf{x}; \boldsymbol{\theta} + \epsilon\mathbf{g})$  to burn in, starting from samples from  $p(\mathbf{x}; \boldsymbol{\theta})$ . Perhaps 1 for RBM on a small image patch, or 5-50 for a more complicated model like a DBM. Initialize a set of  $m$  samples  $\{\tilde{\mathbf{x}}^{(1)}, \dots, \tilde{\mathbf{x}}^{(m)}\}$  to random values (e.g., from a uniform or normal distribution, or possibly a distribution with marginals matched to the model's marginals).

**while** not converged **do**

    Sample a minibatch of  $m$  examples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  from the training set.

$$\mathbf{g} \leftarrow \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}^{(i)}; \boldsymbol{\theta}).$$

**for**  $i = 1$  to  $k$  **do**

**for**  $j = 1$  to  $m$  **do**

$$\tilde{\mathbf{x}}^{(j)} \leftarrow \text{gibbs\_update}(\tilde{\mathbf{x}}^{(j)}).$$

**end for**

**end for**

$$\mathbf{g} \leftarrow \mathbf{g} - \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\tilde{\mathbf{x}}^{(i)}; \boldsymbol{\theta}).$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \epsilon\mathbf{g}.$$

**end while**

---

Care must be taken when evaluating the samples from a model trained with SML. It is necessary to draw the samples starting from a fresh Markov chain initialized from a random starting point after the model is done training. The samples present in the persistent negative chains used for training have been influenced by several recent versions of the model, and thus can make the model appear to have greater capacity than it actually does.

Berglund and Raiko (2013) performed experiments to examine the bias and variance in the estimate of the gradient provided by CD and SML. CD proves to have lower variance than the estimator based on exact sampling. SML has higher variance. The cause of CD’s low variance is its use of the same training points in both the positive and negative phase. If the negative phase is initialized from different training points, the variance rises above that of the estimator based on exact sampling.

All of these methods based on using MCMC to draw samples from the model can in principle be used with almost any variant of MCMC. This means that techniques such as SML can be improved by using any of the enhanced MCMC techniques described in Chapter 17, such as parallel tempering (Desjardins *et al.*, 2010; Cho *et al.*, 2010).

One approach to accelerating mixing during learning relies not on changing the Monte Carlo sampling technology but rather on changing the parametrization of the model and the cost function. *Fast PCD* or *FPCD* (Tieleman and Hinton, 2009) involves replacing the parameters  $\theta$  of a traditional model with an expression

$$\theta = \theta^{(\text{slow})} + \theta^{(\text{fast})}. \quad (18.16)$$

There are now twice as many parameters as before, and they are added together element-wise to provide the parameters used by the original model definition. The fast copy of the parameters is trained with a much larger learning rate, allowing it to adapt rapidly in response to the negative phase of learning and push the Markov chain to new territory. This forces the Markov chain to mix rapidly, though this effect only occurs during learning while the fast weights are free to change. Typically one also applies significant weight decay to the fast weights, encouraging them to converge to small values, after only transiently taking on large values long enough to encourage the Markov chain to change modes.

One key benefit to the MCMC-based methods described in this section is that they provide an estimate of the gradient of  $\log Z$ , and thus we can essentially decompose the problem into the  $\log \tilde{p}$  contribution and the  $\log Z$  contribution. We can then use any other method to tackle  $\log \tilde{p}(\mathbf{x})$ , and just add our negative phase gradient onto the other method’s gradient. In particular, this means that

our positive phase can make use of methods that provide only a lower bound on  $\tilde{p}$ . Most of the other methods of dealing with  $\log Z$  presented in this chapter are incompatible with bound-based positive phase methods.

### 18.3 Pseudolikelihood

Monte Carlo approximations to the partition function and its gradient directly confront the partition function. Other approaches sidestep the issue, by training the model without computing the partition function. Most of these approaches are based on the observation that it is easy to compute ratios of probabilities in an undirected probabilistic model. This is because the partition function appears in both the numerator and the denominator of the ratio and cancels out:

$$\frac{p(\mathbf{x})}{p(\mathbf{y})} = \frac{\frac{1}{Z}\tilde{p}(\mathbf{x})}{\frac{1}{Z}\tilde{p}(\mathbf{y})} = \frac{\tilde{p}(\mathbf{x})}{\tilde{p}(\mathbf{y})}. \quad (18.17)$$

The pseudolikelihood is based on the observation that conditional probabilities take this ratio-based form, and thus can be computed without knowledge of the partition function. Suppose that we partition  $\mathbf{x}$  into  $\mathbf{a}$ ,  $\mathbf{b}$  and  $\mathbf{c}$ , where  $\mathbf{a}$  contains the variables we want to find the conditional distribution over,  $\mathbf{b}$  contains the variables we want to condition on, and  $\mathbf{c}$  contains the variables that are not part of our query.

$$p(\mathbf{a} \mid \mathbf{b}) = \frac{p(\mathbf{a}, \mathbf{b})}{p(\mathbf{b})} = \frac{p(\mathbf{a}, \mathbf{b})}{\sum_{\mathbf{a}, \mathbf{c}} p(\mathbf{a}, \mathbf{b}, \mathbf{c})} = \frac{\tilde{p}(\mathbf{a}, \mathbf{b})}{\sum_{\mathbf{a}, \mathbf{c}} \tilde{p}(\mathbf{a}, \mathbf{b}, \mathbf{c})}. \quad (18.18)$$

This quantity requires marginalizing out  $\mathbf{a}$ , which can be a very efficient operation provided that  $\mathbf{a}$  and  $\mathbf{c}$  do not contain very many variables. In the extreme case,  $\mathbf{a}$  can be a single variable and  $\mathbf{c}$  can be empty, making this operation require only as many evaluations of  $\tilde{p}$  as there are values of a single random variable.

Unfortunately, in order to compute the log-likelihood, we need to marginalize out large sets of variables. If there are  $n$  variables total, we must marginalize a set of size  $n - 1$ . By the chain rule of probability,

$$\log p(\mathbf{x}) = \log p(x_1) + \log p(x_2 \mid x_1) + \cdots + p(x_n \mid \mathbf{x}_{1:n-1}). \quad (18.19)$$

In this case, we have made  $\mathbf{a}$  maximally small, but  $\mathbf{c}$  can be as large as  $\mathbf{x}_{2:n}$ . What if we simply move  $\mathbf{c}$  into  $\mathbf{b}$  to reduce the computational cost? This yields the *pseudolikelihood* (Besag, 1975) objective function, based on predicting the value of

feature  $x_i$  given all of the other features  $\mathbf{x}_{-i}$ :

$$\sum_{i=1}^n \log p(x_i \mid \mathbf{x}_{-i}). \quad (18.20)$$

If each random variable has  $k$  different values, this requires only  $k \times n$  evaluations of  $\tilde{p}$  to compute, as opposed to the  $k^n$  evaluations needed to compute the partition function.

This may look like an unprincipled hack, but it can be proven that estimation by maximizing the pseudolikelihood is asymptotically consistent (Mase, 1995). Of course, in the case of datasets that do not approach the large sample limit, pseudolikelihood may display different behavior from the maximum likelihood estimator.

It is possible to trade computational complexity for deviation from maximum likelihood behavior by using the *generalized pseudolikelihood* estimator (Huang and Ogata, 2002). The generalized pseudolikelihood estimator uses  $m$  different sets  $\mathbb{S}^{(i)}, i = 1, \dots, m$  of indices of variables that appear together on the left side of the conditioning bar. In the extreme case of  $m = 1$  and  $\mathbb{S}^{(1)} = 1, \dots, n$  the generalized pseudolikelihood recovers the log-likelihood. In the extreme case of  $m = n$  and  $\mathbb{S}^{(i)} = \{i\}$ , the generalized pseudolikelihood recovers the pseudolikelihood. The generalized pseudolikelihood objective function is given by

$$\sum_{i=1}^m \log p(\mathbf{x}_{\mathbb{S}^{(i)}} \mid \mathbf{x}_{-\mathbb{S}^{(i)}}). \quad (18.21)$$

The performance of pseudolikelihood-based approaches depends largely on how the model will be used. Pseudolikelihood tends to perform poorly on tasks that require a good model of the full joint  $p(\mathbf{x})$ , such as density estimation and sampling. However, it can perform better than maximum likelihood for tasks that require only the conditional distributions used during training, such as filling in small amounts of missing values. Generalized pseudolikelihood techniques are especially powerful if the data has regular structure that allows the  $\mathbb{S}$  index sets to be designed to capture the most important correlations while leaving out groups of variables that only have negligible correlation. For example, in natural images, pixels that are widely separated in space also have weak correlation, so the generalized pseudolikelihood can be applied with each  $\mathbb{S}$  set being a small, spatially localized window.

One weakness of the pseudolikelihood estimator is that it cannot be used with other approximations that provide only a lower bound on  $\tilde{p}(\mathbf{x})$ , such as variational inference, which will be covered in Chapter 19. This is because  $\tilde{p}$  appears in the

denominator. A lower bound on the denominator provides only an upper bound on the expression as a whole, and there is no benefit to maximizing an upper bound. This makes it difficult to apply pseudolikelihood approaches to deep models such as deep Boltzmann machines, since variational methods are one of the dominant approaches to approximately marginalizing out the many layers of hidden variables that interact with each other. However, pseudolikelihood is still useful for deep learning, because it can be used to train single layer models, or deep models using approximate inference methods that are not based on lower bounds.

Pseudolikelihood has a much greater cost per gradient step than SML, due to its explicit computation of all of the conditionals. However, generalized pseudolikelihood and similar criteria can still perform well if only one randomly selected conditional is computed per example (Goodfellow *et al.*, 2013b), thereby bringing the computational cost down to match that of SML.

Though the pseudolikelihood estimator does not explicitly minimize  $\log Z$ , it can still be thought of as having something resembling a negative phase. The denominators of each conditional distribution result in the learning algorithm suppressing the probability of all states that have only one variable differing from a training example.

See Marlin and de Freitas (2011) for a theoretical analysis of the asymptotic efficiency of pseudolikelihood.

## 18.4 Score Matching and Ratio Matching

Score matching (Hyvärinen, 2005) provides another consistent means of training a model without estimating  $Z$  or its derivatives. The name score matching comes from terminology in which the derivatives of a log density with respect to its argument,  $\nabla_{\mathbf{x}} \log p(\mathbf{x})$ , are called its *score*. The strategy used by score matching is to minimize the expected squared difference between the derivatives of the model's log density with respect to the input and the derivatives of the data's log density with respect to the input:

$$L(\mathbf{x}, \boldsymbol{\theta}) = \frac{1}{2} \|\nabla_{\mathbf{x}} \log p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}) - \nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})\|_2^2 \quad (18.22)$$

$$J(\boldsymbol{\theta}) = \frac{1}{2} \mathbb{E}_{p_{\text{data}}(\mathbf{x})} L(\mathbf{x}, \boldsymbol{\theta}) \quad (18.23)$$

$$\boldsymbol{\theta}^* = \min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \quad (18.24)$$

This objective function avoids the difficulties associated with differentiating the partition function  $Z$  because  $Z$  is not a function of  $\mathbf{x}$  and therefore  $\nabla_{\mathbf{x}} Z = 0$ .

Initially, score matching appears to have a new difficulty: computing the score of the data distribution requires knowledge of the true distribution generating the training data,  $p_{\text{data}}$ . Fortunately, minimizing the expected value of  $L(\mathbf{x}, \boldsymbol{\theta})$  is equivalent to minimizing the expected value of

$$\tilde{L}(\mathbf{x}, \boldsymbol{\theta}) = \sum_{j=1}^n \left[ \frac{\partial^2}{\partial x_j^2} \log p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}) + \frac{1}{2} \left( \frac{\partial}{\partial x_j} \log p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}) \right)^2 \right] \quad (18.25)$$

where  $n$  is the dimensionality of  $\mathbf{x}$ .

Because score matching requires taking derivatives with respect to  $\mathbf{x}$ , it is not applicable to models of discrete data. However, the latent variables in the model may be discrete.

Like the pseudolikelihood, score matching only works when we are able to evaluate  $\log \tilde{p}(\mathbf{x})$  and its derivatives directly. It is not compatible with methods that only provide a lower bound on  $\log \tilde{p}(\mathbf{x})$ , because score matching requires the derivatives and second derivatives of  $\log \tilde{p}(\mathbf{x})$  and a lower bound conveys no information about its derivatives. This means that score matching cannot be applied to estimating models with complicated interactions between the hidden units, such as sparse coding models or deep Boltzmann machines. While score matching can be used to pretrain the first hidden layer of a larger model, it has not been applied as a pretraining strategy for the deeper layers of a larger model. This is probably because the hidden layers of such models usually contain some discrete variables.

While score matching does not explicitly have a negative phase, it can be viewed as a version of contrastive divergence using a specific kind of Markov chain ([Hyvärinen, 2007a](#)). The Markov chain in this case is not Gibbs sampling, but rather a different approach that makes local moves guided by the gradient. Score matching is equivalent to CD with this type of Markov chain when the size of the local moves approaches zero.

[Lyu \(2009\)](#) generalized score matching to the discrete case (but made an error in their derivation that was corrected by [Marlin et al. \(2010\)](#)). [Marlin et al. \(2010\)](#) found that *generalized score matching* (GSM) does not work in high dimensional discrete spaces where the observed probability of many events is 0.

A more successful approach to extending the basic ideas of score matching to discrete data is *ratio matching* ([Hyvärinen, 2007b](#)). Ratio matching applies specifically to binary data. Ratio matching consists of minimizing the average over

examples of the following objective function:

$$L^{(\text{RM})}(\boldsymbol{x}, \boldsymbol{\theta}) = \sum_{j=1}^n \left( \frac{1}{1 + \frac{p_{\text{model}}(\boldsymbol{x}; \boldsymbol{\theta})}{p_{\text{model}}(f(\boldsymbol{x}), j); \boldsymbol{\theta})}} \right)^2, \quad (18.26)$$

where  $f(\boldsymbol{x}, j)$  returns  $\boldsymbol{x}$  with the bit at position  $j$  flipped. Ratio matching avoids the partition function using the same trick as the pseudolikelihood estimator: in a ratio of two probabilities, the partition function cancels out. Marlin *et al.* (2010) found that ratio matching outperforms SML, pseudolikelihood and GSM in terms of the ability of models trained with ratio matching to denoise test set images.

Like the pseudolikelihood estimator, ratio matching requires  $n$  evaluations of  $\tilde{p}$  per data point, making its computational cost per update roughly  $n$  times higher than that of SML.

As with the pseudolikelihood estimator, ratio matching can be thought of as pushing down on all fantasy states that have only one variable different from a training example. Since ratio matching applies specifically to binary data, this means that it acts on all fantasy states within Hamming distance 1 of the data.

Ratio matching can also be useful as the basis for dealing with high-dimensional sparse data, such as word count vectors. This kind of data poses a challenge for MCMC-based methods because the data is extremely expensive to represent in dense format, yet the MCMC sampler does not yield sparse values until the model has learned to represent the sparsity in the data distribution. Dauphin and Bengio (2013) overcame this issue by designing an unbiased stochastic approximation to ratio matching. The approximation evaluates only a randomly selected subset of the terms of the objective, and does not require the model to generate complete fantasy samples.

See Marlin and de Freitas (2011) for a theoretical analysis of the asymptotic efficiency of ratio matching.

## 18.5 Denoising Score Matching

In some cases we may wish to regularize score matching, by fitting a distribution

$$p_{\text{smoothed}}(\boldsymbol{x}) = \int p_{\text{data}}(\boldsymbol{y})q(\boldsymbol{x} \mid \boldsymbol{y})d\boldsymbol{y} \quad (18.27)$$

rather than the true  $p_{\text{data}}$ . The distribution  $q(\boldsymbol{x} \mid \boldsymbol{y})$  is a corruption process, usually one that forms  $\boldsymbol{x}$  by adding a small amount of noise to  $\boldsymbol{y}$ .

Denoising score matching is especially useful because in practice we usually do not have access to the true  $p_{\text{data}}$  but rather only an empirical distribution defined by samples from it. Any consistent estimator will, given enough capacity, make  $p_{\text{model}}$  into a set of Dirac distributions centered on the training points. Smoothing by  $q$  helps to reduce this problem, at the loss of the asymptotic consistency property described in Sec. 5.4.5. Kingma and LeCun (2010) introduced a procedure for performing regularized score matching with the smoothing distribution  $q$  being normally distributed noise.

Recall from Sec. 14.5.1 that several autoencoder training algorithms are equivalent to score matching or denoising score matching. These autoencoder training algorithms are therefore a way of overcoming the partition function problem.

## 18.6 Noise-Contrastive Estimation

Most techniques for estimating models with intractable partition functions do not provide an estimate of the partition function. SML and CD estimate only the gradient of the log partition function, rather than the partition function itself. Score matching and pseudolikelihood avoid computing quantities related to the partition function altogether.

*Noise-contrastive estimation (NCE)* (Gutmann and Hyvärinen, 2010) takes a different strategy. In this approach, the probability distribution estimated by the model is represented explicitly as

$$\log p_{\text{model}}(\mathbf{x}) = \log \tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}) + c, \quad (18.28)$$

where  $c$  is explicitly introduced as an approximation of  $-\log Z(\boldsymbol{\theta})$ . Rather than estimating only  $\boldsymbol{\theta}$ , the noise contrastive estimation procedure treats  $c$  as just another parameter and estimates  $\boldsymbol{\theta}$  and  $c$  simultaneously, using the same algorithm for both. The resulting  $\log p_{\text{model}}(\mathbf{x})$  thus may not correspond exactly to a valid probability distribution, but will become closer and closer to being valid as the estimate of  $c$  improves.<sup>1</sup>

Such an approach would not be possible using maximum likelihood as the criterion for the estimator. The maximum likelihood criterion would choose to set  $c$  arbitrarily high, rather than setting  $c$  to create a valid probability distribution.

---

<sup>1</sup>NCE is also applicable to problems with a tractable partition function, where there is no need to introduce the extra parameter  $c$ . However, it has generated the most interest as a means of estimating models with difficult partition functions.

NCE works by reducing the unsupervised learning problem of estimating  $p(\mathbf{x})$  to that of learning a probabilistic binary classifier in which one of the categories corresponds to the data generated by the model. This supervised learning problem is constructed in such a way that maximum likelihood estimation in this supervised learning problem defines an asymptotically consistent estimator of the original problem.

Specifically, we introduce a second distribution, the *noise distribution*  $p_{\text{noise}}(\mathbf{x})$ . The noise distribution should be tractable to evaluate and to sample from. We can now construct a model over both  $\mathbf{x}$  and a new, binary class variable  $y$ . In the new joint model, we specify that

$$p_{\text{joint}}(y = 1) = \frac{1}{2}, \quad (18.29)$$

$$p_{\text{joint}}(\mathbf{x} \mid y = 1) = p_{\text{model}}(\mathbf{x}), \quad (18.30)$$

and

$$p_{\text{joint}}(\mathbf{x} \mid y = 0) = p_{\text{noise}}(\mathbf{x}). \quad (18.31)$$

In other words,  $y$  is a switch variable that determines whether we will generate  $\mathbf{x}$  from the model or from the noise distribution.

We can construct a similar joint model of training data. In this case, the switch variable determines whether we draw  $\mathbf{x}$  from the *data* or from the noise distribution. Formally,  $p_{\text{train}}(y = 1) = \frac{1}{2}$ ,  $p_{\text{train}}(\mathbf{x} \mid y = 1) = p_{\text{data}}(\mathbf{x})$ , and  $p_{\text{train}}(\mathbf{x} \mid y = 0) = p_{\text{noise}}(\mathbf{x})$ .

We can now just use standard maximum likelihood learning on the *supervised* learning problem of fitting  $p_{\text{joint}}$  to  $p_{\text{train}}$ :

$$\theta, c = \arg \max_{\theta, c} \mathbb{E}_{\mathbf{x}, y \sim p_{\text{train}}} \log p_{\text{joint}}(y \mid \mathbf{x}). \quad (18.32)$$

The distribution  $p_{\text{joint}}$  is essentially a logistic regression model applied to the difference in log probabilities of the model and the noise distribution:

$$p_{\text{joint}}(y = 1 \mid \mathbf{x}) = \frac{p_{\text{model}}(\mathbf{x})}{p_{\text{model}}(\mathbf{x}) + p_{\text{noise}}(\mathbf{x})} \quad (18.33)$$

$$= \frac{1}{1 + \frac{p_{\text{noise}}(\mathbf{x})}{p_{\text{model}}(\mathbf{x})}} \quad (18.34)$$

$$= \frac{1}{1 + \exp \left( \log \frac{p_{\text{noise}}(\mathbf{x})}{p_{\text{model}}(\mathbf{x})} \right)} \quad (18.35)$$

$$= \sigma \left( -\log \frac{p_{\text{noise}}(\mathbf{x})}{p_{\text{model}}(\mathbf{x})} \right) \quad (18.36)$$

$$= \sigma (\log p_{\text{model}}(\mathbf{x}) - \log p_{\text{noise}}(\mathbf{x})). \quad (18.37)$$

NCE is thus simple to apply so long as  $\log \tilde{p}_{\text{model}}$  is easy to back-propagate through, and, as specified above,  $p_{\text{noise}}$  is easy to evaluate (in order to evaluate  $p_{\text{joint}}$ ) and sample from (in order to generate the training data).

NCE is most successful when applied to problems with few random variables, but can work well even if those random variables can take on a high number of values. For example, it has been successfully applied to modeling the conditional distribution over a word given the context of the word (Mnih and Kavukcuoglu, 2013). Though the word may be drawn from a large vocabulary, there is only one word.

When NCE is applied to problems with many random variables, it becomes less efficient. The logistic regression classifier can reject a noise sample by identifying any one variable whose value is unlikely. This means that learning slows down greatly after  $p_{\text{model}}$  has learned the basic marginal statistics. Imagine learning a model of images of faces, using unstructured Gaussian noise as  $p_{\text{noise}}$ . If  $p_{\text{model}}$  learns about eyes, it can reject almost all unstructured noise samples without having learned anything about other facial features, such as mouths.

The constraint that  $p_{\text{noise}}$  must be easy to evaluate and easy to sample from can be overly restrictive. When  $p_{\text{noise}}$  is simple, most samples are likely to be too obviously distinct from the data to force  $p_{\text{model}}$  to improve noticeably.

Like score matching and pseudolikelihood, NCE does not work if only a lower bound on  $\tilde{p}$  is available. Such a lower bound could be used to construct a lower bound on  $p_{\text{joint}}(y = 1 | \mathbf{x})$ , but it can only be used to construct an upper bound on  $p_{\text{joint}}(y = 0 | \mathbf{x})$ , which appears in half the terms of the NCE objective. Likewise, a lower bound on  $p_{\text{noise}}$  is not useful, because it provides only an upper bound on  $p_{\text{joint}}(y = 1 | \mathbf{x})$ .

When the model distribution is copied to define a new noise distribution before each gradient step, NCE defines a procedure called *self-contrastive estimation*, whose expected gradient is equivalent to the expected gradient of maximum likelihood (Goodfellow, 2014). The special case of NCE where the noise samples are those generated by the model suggests that maximum likelihood can be interpreted as a procedure that forces a model to constantly learn to distinguish reality from its own evolving beliefs, while noise contrastive estimation achieves some reduced computational cost by only forcing the model to distinguish reality from a fixed baseline (the noise model).

Using the supervised task of classifying between training samples and generated samples (with the model energy function used in defining the classifier) to provide a gradient on the model was introduced earlier in various forms (Welling *et al.*, 2003b; Bengio, 2009).

Noise contrastive estimation is based on the idea that a good generative model should be able to distinguish data from noise. A closely related idea is that a good generative model should be able to generate samples that no classifier can distinguish from data. This idea yields generative adversarial networks (Sec. 20.10.4).

## 18.7 Estimating the Partition Function

While much of this chapter is dedicated to describing methods that avoid needing to compute the intractable partition function  $Z(\boldsymbol{\theta})$  associated with an undirected graphical model, in this section we discuss several methods for directly estimating the partition function.

Estimating the partition function can be important because we require it if we wish to compute the normalized likelihood of data. This is often important in *evaluating* the model, monitoring training performance, and comparing models to each other.

For example, imagine we have two models: model  $\mathcal{M}_A$  defining a probability distribution  $p_A(\mathbf{x}; \boldsymbol{\theta}_A) = \frac{1}{Z_A} \tilde{p}_A(\mathbf{x}; \boldsymbol{\theta}_A)$  and model  $\mathcal{M}_B$  defining a probability distribution  $p_B(\mathbf{x}; \boldsymbol{\theta}_B) = \frac{1}{Z_B} \tilde{p}_B(\mathbf{x}; \boldsymbol{\theta}_B)$ . A common way to compare the models is to evaluate and compare the likelihood that both models assign to an i.i.d. test dataset. Suppose the test set consists of  $m$  examples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ . If  $\prod_i p_A(\mathbf{x}^{(i)}; \boldsymbol{\theta}_A) > \prod_i p_B(\mathbf{x}^{(i)}; \boldsymbol{\theta}_B)$  or equivalently if

$$\sum_i \log p_A(\mathbf{x}^{(i)}; \boldsymbol{\theta}_A) - \sum_i \log p_B(\mathbf{x}^{(i)}; \boldsymbol{\theta}_B) > 0, \quad (18.38)$$

then we say that  $\mathcal{M}_A$  is a better model than  $\mathcal{M}_B$  (or, at least, it is a better model of the test set), in the sense that it has a better test log-likelihood. Unfortunately, testing whether this condition holds requires knowledge of the partition function. Unfortunately, Eq. 18.38 seems to require evaluating the log probability that the model assigns to each point, which in turn requires evaluating the partition function. We can simplify the situation slightly by re-arranging Eq. 18.38 into a

form where we need to know only the **ratio** of the two model's partition functions:

$$\sum_i \log p_A(\mathbf{x}^{(i)}; \boldsymbol{\theta}_A) - \sum_i \log p_B(\mathbf{x}^{(i)}; \boldsymbol{\theta}_B) = \sum_i \left( \log \frac{\tilde{p}_A(\mathbf{x}^{(i)}; \boldsymbol{\theta}_A)}{\tilde{p}_B(\mathbf{x}^{(i)}; \boldsymbol{\theta}_B)} \right) - m \log \frac{Z(\boldsymbol{\theta}_A)}{Z(\boldsymbol{\theta}_B)}. \quad (18.39)$$

We can thus determine whether  $\mathcal{M}_A$  is a better model than  $\mathcal{M}_B$  without knowing the partition function of either model but only their ratio. As we will see shortly, we can estimate this ratio using importance sampling, provided that the two models are similar.

If, however, we wanted to compute the actual probability of the test data under either  $\mathcal{M}_A$  or  $\mathcal{M}_B$ , we would need to compute the actual value of the partition functions. That said, if we knew the ratio of two partition functions,  $r = \frac{Z(\boldsymbol{\theta}_B)}{Z(\boldsymbol{\theta}_A)}$ , and we knew the actual value of just one of the two, say  $Z(\boldsymbol{\theta}_A)$ , we could compute the value of the other:

$$Z(\boldsymbol{\theta}_B) = rZ(\boldsymbol{\theta}_A) = \frac{Z(\boldsymbol{\theta}_B)}{Z(\boldsymbol{\theta}_A)} Z(\boldsymbol{\theta}_A). \quad (18.40)$$

A simple way to estimate the partition function is to use a Monte Carlo method such as simple importance sampling. We present the approach in terms of continuous variables using integrals, but it can be readily applied to discrete variables by replacing the integrals with summation. We use a proposal distribution  $p_0(\mathbf{x}) = \frac{1}{Z_0} \tilde{p}_0(\mathbf{x})$  which supports tractable sampling and tractable evaluation of both the partition function  $Z_0$  and the unnormalized distribution  $\tilde{p}_0(\mathbf{x})$ .

$$Z_1 = \int \tilde{p}_1(\mathbf{x}) d\mathbf{x} \quad (18.41)$$

$$= \int \frac{p_0(\mathbf{x})}{\tilde{p}_0(\mathbf{x})} \tilde{p}_1(\mathbf{x}) d\mathbf{x} \quad (18.42)$$

$$= Z_0 \int p_0(\mathbf{x}) \frac{\tilde{p}_1(\mathbf{x})}{\tilde{p}_0(\mathbf{x})} d\mathbf{x} \quad (18.43)$$

$$\hat{Z}_1 = \frac{Z_0}{K} \sum_{k=1}^K \frac{\tilde{p}_1(\mathbf{x}^{(k)})}{\tilde{p}_0(\mathbf{x}^{(k)})} \quad \text{s.t. : } \mathbf{x}^{(k)} \sim p_0 \quad (18.44)$$

In the last line, we make a Monte Carlo estimator,  $\hat{Z}_1$ , of the integral using samples drawn from  $p_0(\mathbf{x})$  and then weight each sample with the ratio of the unnormalized  $\tilde{p}_1$  and the proposal  $p_0$ .

We see also that this approach allows us to estimate the ratio between the

partition functions as

$$\frac{1}{K} \sum_{k=1}^K \frac{\tilde{p}_1(\mathbf{x}^{(k)})}{\tilde{p}_0(\mathbf{x}^{(k)})} \quad \text{s.t. : } \mathbf{x}^{(k)} \sim p_0. \quad (18.45)$$

This value can then be used directly to compare two models as described in Eq. 18.39.

If the distribution  $p_0$  is close to  $p_1$ , Eq. 18.44 can be an effective way of estimating the partition function (Minka, 2005). Unfortunately, most of the time  $p_1$  is both complicated (usually multimodal) and defined over a high dimensional space. It is difficult to find a tractable  $p_0$  that is simple enough to evaluate while still being close enough to  $p_1$  to result in a high quality approximation. If  $p_0$  and  $p_1$  are not close, most samples from  $p_0$  will have low probability under  $p_1$  and therefore make (relatively) negligible contribution to the sum in Eq. 18.44.

Having few samples with significant weights in this sum will result in an estimator that is of poor quality due to high variance. This can be understood quantitatively through an estimate of the variance of our estimate  $\hat{Z}_1$ :

$$\text{Var}(\hat{Z}_1) = \frac{Z_0}{K^2} \sum_{k=1}^K \left( \frac{\tilde{p}_1(\mathbf{x}^{(k)})}{\tilde{p}_0(\mathbf{x}^{(k)})} - \hat{Z}_1 \right)^2. \quad (18.46)$$

This quantity is largest when there is significant deviation in the values of the importance weights  $\frac{\tilde{p}_1(\mathbf{x}^{(k)})}{\tilde{p}_0(\mathbf{x}^{(k)})}$ .

We now turn to two related strategies developed to cope with the challenging task of estimating partition functions for complex distributions over high-dimensional spaces: annealed importance sampling and bridge sampling. Both start with the simple importance sampling strategy introduced above and both attempt to overcome the problem of the proposal  $p_0$  being too far from  $p_1$  by introducing intermediate distributions that attempt to *bridge the gap* between  $p_0$  and  $p_1$ .

### 18.7.1 Annealed Importance Sampling

In situations where  $D_{\text{KL}}(p_0 \| p_1)$  is large (i.e., where there is little overlap between  $p_0$  and  $p_1$ ), a strategy called *annealed importance sampling* (AIS) attempts to bridge the gap by introducing intermediate distributions (Jarzynski, 1997; Neal, 2001). Consider a sequence of distributions  $p_{\eta_0}, \dots, p_{\eta_n}$ , with  $0 = \eta_0 < \eta_1 < \dots < \eta_{n-1} < \eta_n = 1$  so that the first and last distributions in the sequence are  $p_0$  and  $p_1$  respectively.

This approach allows us to estimate the partition function of a multimodal distribution defined over a high-dimensional space (such as the distribution defined by a trained RBM). We begin with a simpler model with a known partition function (such as an RBM with zeroes for weights) and estimate the ratio between the two model's partition functions. The estimate of this ratio is based on the estimate of the ratios of a sequence of many similar distributions, such as the sequence of RBMs with weights interpolating between zero and the learned weights.

We can now write the ratio  $\frac{Z_1}{Z_0}$  as

$$\frac{Z_1}{Z_0} = \frac{Z_1}{Z_0} \frac{Z_{\eta_1}}{Z_{\eta_1}} \dots \frac{Z_{\eta_{n-1}}}{Z_{\eta_{n-1}}} \quad (18.47)$$

$$= \frac{Z_{\eta_1}}{Z_0} \frac{Z_{\eta_2}}{Z_{\eta_1}} \dots \frac{Z_{\eta_{n-1}}}{Z_{\eta_{n-2}}} \frac{Z_1}{Z_{\eta_{n-1}}} \quad (18.48)$$

$$= \prod_{j=0}^{n-1} \frac{Z_{\eta_{j+1}}}{Z_{\eta_j}} \quad (18.49)$$

Provided the distributions  $p_{\eta_j}$  and  $p_{\eta_{j+1}}$ , for all  $0 \leq j \leq n - 1$ , are sufficiently close, we can reliably estimate each of the factors  $\frac{Z_{\eta_{j+1}}}{Z_{\eta_j}}$  using simple importance sampling and then use these to obtain an estimate of  $\frac{Z_1}{Z_0}$ .

Where do these intermediate distributions come from? Just as the original proposal distribution  $p_0$  is a design choice, so is the sequence of distributions  $p_{\eta_1} \dots p_{\eta_{n-1}}$ . That is, it can be specifically constructed to suit the problem domain. One general-purpose and popular choice for the intermediate distributions is to use the weighted geometric average of the target distribution  $p_1$  and the starting proposal distribution (for which the partition function is known)  $p_0$ :

$$p_{\eta_j} \propto p_1^{\eta_j} p_0^{1-\eta_j} \quad (18.50)$$

In order to sample from these intermediate distributions, we define a series of Markov chain transition functions  $T_{\eta_j}(\mathbf{x}' | \mathbf{x})$  that define the conditional probability distribution of transitioning to  $\mathbf{x}'$  given we are currently at  $\mathbf{x}$ . The transition operator  $T_{\eta_j}(\mathbf{x}' | \mathbf{x})$  is defined to leave  $p_{\eta_j}(\mathbf{x})$  invariant:

$$p_{\eta_j}(\mathbf{x}) = \int p_{\eta_j}(\mathbf{x}') T_{\eta_j}(\mathbf{x} | \mathbf{x}') d\mathbf{x}' \quad (18.51)$$

These transitions may be constructed as any Markov chain Monte Carlo method (e.g., Metropolis-Hastings, Gibbs), including methods involving multiple passes through all of the random variables or other kinds of iterations.

The AIS sampling strategy is then to generate samples from  $p_0$  and then use the transition operators to sequentially generate samples from the intermediate distributions until we arrive at samples from the target distribution  $p_1$ :

- for  $k = 1 \dots K$ 
  - Sample  $\mathbf{x}_{\eta_1}^{(k)} \sim p_0(\mathbf{x})$
  - Sample  $\mathbf{x}_{\eta_2}^{(k)} \sim T_{\eta_1}(\mathbf{x}_{\eta_2}^{(k)} | \mathbf{x}_{\eta_1}^{(k)})$
  - ...
  - Sample  $\mathbf{x}_{\eta_{n-1}}^{(k)} \sim T_{\eta_{n-2}}(\mathbf{x}_{\eta_{n-1}}^{(k)} | \mathbf{x}_{\eta_{n-2}}^{(k)})$
  - Sample  $\mathbf{x}_{\eta_n}^{(k)} \sim T_{\eta_{n-1}}(\mathbf{x}_{\eta_n}^{(k)} | \mathbf{x}_{\eta_{n-1}}^{(k)})$
- end

For sample  $k$ , we can derive the importance weight by chaining together the importance weights for the jumps between the intermediate distributions given in Eq. 18.49.

$$w^{(k)} = \frac{\tilde{p}_{\eta_1}(\mathbf{x}_{\eta_1}^{(k)})}{\tilde{p}_0(\mathbf{x}_0^{(k)})} \frac{\tilde{p}_{\eta_2}(\mathbf{x}_{\eta_2}^{(k)})}{\tilde{p}_{\eta_1}(\mathbf{x}_{\eta_1}^{(k)})} \dots \frac{\tilde{p}_1(\mathbf{x}_1^{(k)})}{\tilde{p}_{\eta_{n-1}}(\mathbf{x}_{\eta_{n-1}}^{(k)})} \quad (18.52)$$

To avoid computational issues such as overflow, it is probably best to do the computation in log space:

$$\log w^{(k)} = \log \tilde{p}_{\eta_1}(\mathbf{x}) - \log \tilde{p}_0(\mathbf{x}) + \dots \quad (18.53)$$

With the sampling procedure thus defined and the importance weights given in Eq. 18.52, the estimate of the ratio of partition functions is given by:

$$\frac{Z_1}{Z_0} \approx \frac{1}{K} \sum_{k=1}^K w^{(k)} \quad (18.54)$$

In order to verify that this procedure defines a valid importance sampling scheme, we can show (Neal, 2001) that the AIS procedure corresponds to simple importance sampling on an extended state space with points sampled over the product space  $[\mathbf{x}_{\eta_1}, \dots, \mathbf{x}_{\eta_{n-1}}, \mathbf{x}_1]$ . To do this, we define the distribution over the extended space as:

$$\tilde{p}(\mathbf{x}_{\eta_1}, \dots, \mathbf{x}_{\eta_{n-1}}, \mathbf{x}_1) \quad (18.55)$$

$$= \tilde{p}_1(\mathbf{x}_1) \tilde{T}_{\eta_{n-1}}(\mathbf{x}_{\eta_{n-1}} | \mathbf{x}_1) \tilde{T}_{\eta_{n-2}}(\mathbf{x}_{\eta_{n-2}} | \mathbf{x}_{\eta_{n-1}}) \dots \tilde{T}_{\eta_1}(\mathbf{x}_{\eta_1} | \mathbf{x}_{\eta_2}), \quad (18.56)$$

where  $\tilde{T}_a$  is the reverse of the transition operator defined by  $T_a$  (via an application of Bayes' rule):

$$\tilde{T}_a(\mathbf{x}' \mid \mathbf{x}) = \frac{p_a(\mathbf{x}')}{p_a(\mathbf{x})} T_a(\mathbf{x} \mid \mathbf{x}') = \frac{\tilde{p}_a(\mathbf{x}')}{\tilde{p}_a(\mathbf{x})} T_a(\mathbf{x} \mid \mathbf{x}'). \quad (18.57)$$

Plugging the above into the expression for the joint distribution on the extended state space given in Eq. 18.56, we get:

$$\tilde{p}(\mathbf{x}_{\eta_1}, \dots, \mathbf{x}_{\eta_{n-1}}, \mathbf{x}_1) \quad (18.58)$$

$$= \tilde{p}_1(\mathbf{x}_1) \frac{\tilde{p}_{\eta_{n-1}}(\mathbf{x}_{\eta_{n-1}})}{\tilde{p}_{\eta_{n-1}}(\mathbf{x}_1)} T_{\eta_{n-1}}(\mathbf{x}_1 \mid \mathbf{x}_{\eta_{n-1}}) \prod_{i=1}^{n-2} \frac{\tilde{p}_{\eta_i}(\mathbf{x}_{\eta_i})}{\tilde{p}_{\eta_i}(\mathbf{x}_{\eta_{i+1}})} T_{\eta_i}(\mathbf{x}_{\eta_{i+1}} \mid \mathbf{x}_{\eta_i}) \quad (18.59)$$

$$= \frac{\tilde{p}_1(\mathbf{x}_1)}{\tilde{p}_{\eta_{n-1}}(\mathbf{x}_1)} T_{\eta_{n-1}}(\mathbf{x}_1 \mid \mathbf{x}_{\eta_{n-1}}) \tilde{p}_{\eta_1}(\mathbf{x}_{\eta_1}) \prod_{i=1}^{n-2} \frac{\tilde{p}_{\eta_{i+1}}(\mathbf{x}_{\eta_{i+1}})}{\tilde{p}_{\eta_i}(\mathbf{x}_{\eta_{i+1}})} T_{\eta_i}(\mathbf{x}_{\eta_{i+1}} \mid \mathbf{x}_{\eta_i}). \quad (18.60)$$

We now have means of generating samples from the joint proposal distribution  $q$  over the extended sample via a sampling scheme given above, with the joint distribution given by:

$$q(\mathbf{x}_{\eta_1}, \dots, \mathbf{x}_{\eta_{n-1}}, \mathbf{x}_1) = p_0(\mathbf{x}_{\eta_1}) T_{\eta_1}(\mathbf{x}_{\eta_2} \mid \mathbf{x}_{\eta_1}) \dots T_{\eta_{n-1}}(\mathbf{x}_1 \mid \mathbf{x}_{\eta_{n-1}}). \quad (18.61)$$

We have a joint distribution on the extended space given by Eq. 18.60. Taking  $q(\mathbf{x}_{\eta_1}, \dots, \mathbf{x}_{\eta_{n-1}}, \mathbf{x}_1)$  as the proposal distribution on the extended state space from which we will draw samples, it remains to determine the importance weights:

$$w^{(k)} = \frac{\tilde{p}(\mathbf{x}_{\eta_1}, \dots, \mathbf{x}_{\eta_{n-1}}, \mathbf{x}_1)}{q(\mathbf{x}_{\eta_1}, \dots, \mathbf{x}_{\eta_{n-1}}, \mathbf{x}_1)} = \frac{\tilde{p}_1(\mathbf{x}_1^{(k)})}{\tilde{p}_{\eta_{n-1}}(\mathbf{x}_{\eta_{n-1}}^{(k)})} \dots \frac{\tilde{p}_{\eta_2}(\mathbf{x}_{\eta_2}^{(k)})}{\tilde{p}_1(\mathbf{x}_{\eta_1}^{(k)})} \frac{\tilde{p}_{\eta_1}(\mathbf{x}_{\eta_1}^{(k)})}{\tilde{p}_0(\mathbf{x}_0^{(k)})}. \quad (18.62)$$

These weights are the same as proposed for AIS. Thus we can interpret AIS as simple importance sampling applied to an extended state and its validity follows immediately from the validity of importance sampling.

Annealed importance sampling (AIS) was first discovered by [Jarzynski \(1997\)](#) and then again, independently, by [Neal \(2001\)](#). It is currently the most common way of estimating the partition function for undirected probabilistic models. The reasons for this may have more to do with the publication of an influential paper ([Salakhutdinov and Murray, 2008](#)) describing its application to estimating the partition function of restricted Boltzmann machines and deep belief networks than

with any inherent advantage the method has over the other method described below.

A discussion of the properties of the AIS estimator (e.g.. its variance and efficiency) can be found in [Neal \(2001\)](#).

### 18.7.2 Bridge Sampling

Bridge sampling [Bennett \(1976\)](#) is another method that, like AIS, addresses the shortcomings of importance sampling. Rather than chaining together a series of intermediate distributions, bridge sampling relies on a single distribution  $p_*$ , known as the bridge, to interpolate between a distribution with known partition function,  $p_0$ , and a distribution  $p_1$  for which we are trying to estimate the partition function  $Z_1$ .

Bridge sampling estimates the ratio  $Z_1/Z_0$  as the ratio of the expected importance weights between  $\tilde{p}_0$  and  $\tilde{p}_*$  and between  $\tilde{p}_1$  and  $\tilde{p}_*$ :

$$\frac{Z_1}{Z_0} \approx \sum_{k=1}^K \frac{\tilde{p}_*(\mathbf{x}_0^{(k)})}{\tilde{p}_0(\mathbf{x}_0^{(k)})} \Bigg/ \sum_{k=1}^K \frac{\tilde{p}_*(\mathbf{x}_1^{(k)})}{\tilde{p}_1(\mathbf{x}_1^{(k)})} \quad (18.63)$$

If the bridge distribution  $p_*$  is chosen carefully to have a large overlap of support with both  $p_0$  and  $p_1$ , then bridge sampling can allow the distance between two distributions (or more formally,  $D_{\text{KL}}(p_0\|p_1)$ ) to be much larger than with standard importance sampling.

It can be shown that the optimal bridging distribution is given by  $p_*^{(\text{opt})}(\mathbf{x}) \propto \frac{\tilde{p}_0(\mathbf{x})\tilde{p}_1(\mathbf{x})}{r\tilde{p}_0(\mathbf{x}) + \tilde{p}_1(\mathbf{x})}$  where  $r = Z_1/Z_0$ . At first, this appears to be an unworkable solution as it would seem to require the very quantity we are trying to estimate,  $Z_1/Z_0$ . However, it is possible to start with a coarse estimate of  $r$  and use the resulting bridge distribution to refine our estimate iteratively ([Neal, 2005](#)). That is, we iteratively re-estimate the ratio and use each iteration to update the value of  $r$ .

**Linked importance sampling** Both AIS and bridge sampling have their advantages. If  $D_{\text{KL}}(p_0\|p_1)$  is not too large (because  $p_0$  and  $p_1$  are sufficiently close) bridge sampling can be a more effective means of estimating the ratio of partition functions than AIS. If, however, the two distributions are too far apart for a single distribution  $p_*$  to bridge the gap then one can at least use AIS with potentially many intermediate distributions to span the distance between  $p_0$  and  $p_1$ . [Neal \(2005\)](#) showed how his linked importance sampling method leveraged the power of the bridge sampling strategy to bridge the intermediate distributions used in AIS to significantly improve the overall partition function estimates.

**Estimating the partition function while training** While AIS has become accepted as the standard method for estimating the partition function for many undirected models, it is sufficiently computationally intensive that it remains infeasible to use during training. However, alternative strategies that have been explored to maintain an estimate of the partition function throughout training

Using a combination of bridge sampling, short-chain AIS and parallel tempering, Desjardins *et al.* (2011) devised a scheme to track the partition function of an RBM throughout the training process. The strategy is based on the maintenance of independent estimates of the partition functions of the RBM at every temperature operating in the parallel tempering scheme. The authors combined bridge sampling estimates of the ratios of partition functions of neighboring chains (i.e. from parallel tempering) with AIS estimates across time to come up with a low variance estimate of the partition functions at every iteration of learning.

The tools described in this chapter provide many different ways of overcoming the problem of intractable partition functions, but there can be several other difficulties involved in training and using generative models. Foremost among these is the problem of intractable inference, which we confront next.

# Chapter 19

## Approximate inference

Many probabilistic models are difficult to train because it is difficult to perform inference in them. In the context of deep learning, we usually have a set of visible variables  $\mathbf{v}$  and a set of latent variables  $\mathbf{h}$ . The challenge of inference usually refers to the difficult problem of computing  $p(\mathbf{h} \mid \mathbf{v})$  or taking expectations with respect to it. Such operations are often necessary for tasks like maximum likelihood learning.

Many simple graphical models with only one hidden layer, such as restricted Boltzmann machines and probabilistic PCA, are defined in a way that makes inference operations like computing  $p(\mathbf{h} \mid \mathbf{v})$ , or taking expectations with respect to it, simple. Unfortunately, most graphical models with multiple layers of hidden variables have intractable posterior distributions. Exact inference requires an exponential amount of time in these models. Even some models with only a single layer, such as sparse coding, have this problem.

In this chapter, we introduce several of the techniques for confronting these intractable inference problems. Later, in Chapter 20, we will describe how to use these techniques to train probabilistic models that would otherwise be intractable, such as deep belief networks and deep Boltzmann machines.

Intractable inference problems in deep learning usually arise from interactions between latent variables in a structured graphical model. See Fig. 19.1 for some examples. These interactions may be due to direct interactions in undirected models or “explaining away” interactions between mutual ancestors of the same visible unit in directed models.

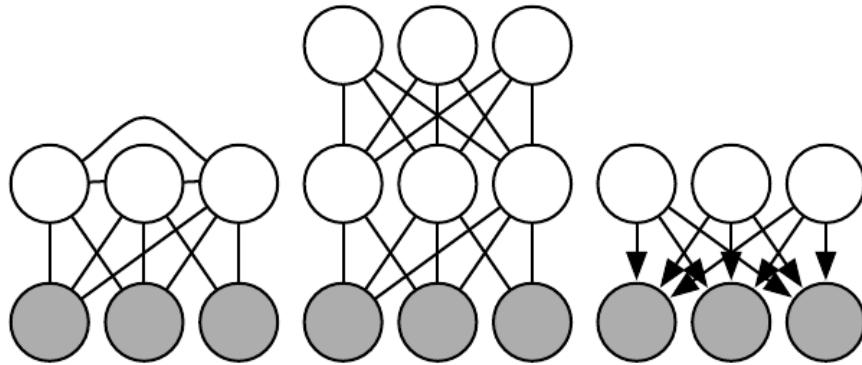


Figure 19.1: Intractable inference problems in deep learning are usually the result of interactions between latent variables in a structured graphical model. These can be due to edges directly connecting one latent variable to another, or due to longer paths that are activated when the child of a V-structure is observed. (*Left*) A *semi-restricted Boltzmann machine* (Osindero and Hinton, 2008) with connections between hidden units. These direct connections between latent variables make the posterior distribution intractable due to large cliques of latent variables. (*Center*) A deep Boltzmann machine, organized into layers of variables without intra-layer connections, still has an intractable posterior distribution due to the connections between layers. (*Right*) This directed model has interactions between latent variables when the visible variables are observed, because every two latent variables are co-parents. Some probabilistic models are able to provide tractable inference over the latent variables despite having one of the graph structures depicted above. This is possible if the conditional probability distributions are chosen to introduce additional independences beyond those described by the graph. For example, probabilistic PCA has the graph structure shown in the right, yet still has simple inference due to special properties of the specific conditional distributions it uses (linear-Gaussian conditionals with mutually orthogonal basis vectors).

## 19.1 Inference as Optimization

Many approaches to confronting the problem of difficult inference make use of the observation that exact inference can be described as an optimization problem. Approximate inference algorithms may then be derived by approximating the underlying optimization problem.

To construct the optimization problem, assume we have a probabilistic model consisting of observed variables  $\mathbf{v}$  and latent variables  $\mathbf{h}$ . We would like to compute the log probability of the observed data,  $\log p(\mathbf{v}; \boldsymbol{\theta})$ . Sometimes it is too difficult to compute  $\log p(\mathbf{v}; \boldsymbol{\theta})$  if it is costly to marginalize out  $\mathbf{h}$ . Instead, we can compute a lower bound  $\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q)$  on  $\log p(\mathbf{v}; \boldsymbol{\theta})$ . This bound is called the *evidence lower bound* (ELBO). Another commonly used name for this lower bound is the negative *variational free energy*. Specifically, the evidence lower bound is defined to be

$$\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q) = \log p(\mathbf{v}; \boldsymbol{\theta}) - D_{\text{KL}}(q(\mathbf{h} | \mathbf{v}) \| p(\mathbf{h} | \mathbf{v}; \boldsymbol{\theta})) \quad (19.1)$$

where  $q$  is an arbitrary probability distribution over  $\mathbf{h}$ .

Because the difference between  $\log p(\mathbf{v})$  and  $\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q)$  is given by the KL divergence and because the KL divergence is always non-negative, we can see that  $\mathcal{L}$  always has at most the same value as the desired log probability. The two are equal if and only if  $q$  is the same distribution as  $p(\mathbf{h} | \mathbf{v})$ .

Surprisingly,  $\mathcal{L}$  can be considerably easier to compute for some distributions  $q$ . Simple algebra shows that we can rearrange  $\mathcal{L}$  into a much more convenient form:

$$\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q) = \log p(\mathbf{v}; \boldsymbol{\theta}) - D_{\text{KL}}(q(\mathbf{h} | \mathbf{v}) \| p(\mathbf{h} | \mathbf{v}; \boldsymbol{\theta})) \quad (19.2)$$

$$= \log p(\mathbf{v}; \boldsymbol{\theta}) - \mathbb{E}_{\mathbf{h} \sim q} \log \frac{q(\mathbf{h} | \mathbf{v})}{p(\mathbf{h} | \mathbf{v})} \quad (19.3)$$

$$= \log p(\mathbf{v}; \boldsymbol{\theta}) - \mathbb{E}_{\mathbf{h} \sim q} \log \frac{q(\mathbf{h} | \mathbf{v})}{\frac{p(\mathbf{h}, \mathbf{v}; \boldsymbol{\theta})}{p(\mathbf{v}; \boldsymbol{\theta})}} \quad (19.4)$$

$$= \log p(\mathbf{v}; \boldsymbol{\theta}) - \mathbb{E}_{\mathbf{h} \sim q} [\log q(\mathbf{h} | \mathbf{v}) - \log p(\mathbf{h}, \mathbf{v}; \boldsymbol{\theta}) + \log p(\mathbf{v}; \boldsymbol{\theta})] \quad (19.5)$$

$$= -\mathbb{E}_{\mathbf{h} \sim q} [\log q(\mathbf{h} | \mathbf{v}) - \log p(\mathbf{h}, \mathbf{v}; \boldsymbol{\theta})]. \quad (19.6)$$

This yields the more canonical definition of the evidence lower bound,

$$\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q) = \mathbb{E}_{\mathbf{h} \sim q} [\log p(\mathbf{h}, \mathbf{v})] + H(q). \quad (19.7)$$

For an appropriate choice of  $q$ ,  $\mathcal{L}$  is tractable to compute. For any choice of  $q$ ,  $\mathcal{L}$  provides a lower bound on the likelihood. For  $q(\mathbf{h} | \mathbf{v})$  that are better

approximations of  $p(\mathbf{h} \mid \mathbf{v})$ , the lower bound  $\mathcal{L}$  will be tighter, in other words, closer to  $\log p(\mathbf{v})$ . When  $q(\mathbf{h} \mid \mathbf{v}) = p(\mathbf{h} \mid \mathbf{v})$ , the approximation is perfect, and  $\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q) = \log p(\mathbf{v}; \boldsymbol{\theta})$ .

We can thus think of inference as the procedure for finding the  $q$  that maximizes  $\mathcal{L}$ . Exact inference maximizes  $\mathcal{L}$  perfectly by searching over a family of functions  $q$  that includes  $p(\mathbf{h} \mid \mathbf{v})$ . Throughout this chapter, we will show how to derive different forms of approximate inference by using approximate optimization to find  $q$ . We can make the optimization procedure less expensive but approximate by restricting the family of distributions  $q$  the optimization is allowed to search over or by using an imperfect optimization procedure that may not completely maximize  $\mathcal{L}$  but merely increase it by a significant amount.

No matter what choice of  $q$  we use,  $\mathcal{L}$  is a lower bound. We can get tighter or looser bounds that are cheaper or more expensive to compute depending on how we choose to approach this optimization problem. We can obtain a poorly matched  $q$  but reduce the computational cost by using an imperfect optimization procedure, or by using a perfect optimization procedure over a restricted family of  $q$  distributions.

## 19.2 Expectation Maximization

The first algorithm we introduce based on maximizing a lower bound  $\mathcal{L}$  is the *expectation maximization* (EM) algorithm, a popular training algorithm for models with latent variables. We describe here a view on the EM algorithm developed by [Neal and Hinton \(1999\)](#). Unlike most of the other algorithms we describe in this chapter, EM is not an approach to approximate inference, but rather an approach to learning with an approximate posterior.

The EM algorithm consists of alternating between two steps until convergence:

- The *E-step* (Expectation step): Let  $\boldsymbol{\theta}^{(0)}$  denote the value of the parameters at the beginning of the step. Set  $q(\mathbf{h}^{(i)} \mid \mathbf{v}) = p(\mathbf{h}^{(i)} \mid \mathbf{v}^{(i)}; \boldsymbol{\theta}^{(0)})$  for all indices  $i$  of the training examples  $\mathbf{v}^{(i)}$  we want to train on (both batch and minibatch variants are valid). By this we mean  $q$  is defined in terms of the **current** parameter value of  $\boldsymbol{\theta}^{(0)}$ ; if we vary  $\boldsymbol{\theta}$  then  $p(\mathbf{h} \mid \mathbf{v}; \boldsymbol{\theta})$  will change but  $q(\mathbf{h} \mid \mathbf{v})$  will remain equal to  $p(\mathbf{h} \mid \mathbf{v}; \boldsymbol{\theta}^{(0)})$ .
- The M-step (Maximization step): Completely or partially maximize

$$\sum_i \mathcal{L}(\mathbf{v}^{(i)}, \boldsymbol{\theta}, q) \tag{19.8}$$

with respect to  $\theta$  using your optimization algorithm of choice.

This can be viewed as a coordinate ascent algorithm to maximize  $\mathcal{L}$ . On one step, we maximize  $\mathcal{L}$  with respect to  $q$ , and on the other, we maximize  $\mathcal{L}$  with respect to  $\theta$ .

Stochastic gradient ascent on latent variable models can be seen as a special case of the EM algorithm where the M step consists of taking a single gradient step. Other variants of the EM algorithm can make much larger steps. For some model families, the M step can even be performed analytically, jumping all the way to the optimal solution for  $\theta$  given the current  $q$ .

Even though the E-step involves exact inference, we can think of the EM algorithm as using approximate inference in some sense. Specifically, the M-step assumes that the same value of  $q$  can be used for all values of  $\theta$ . This will introduce a gap between  $\mathcal{L}$  and the true  $\log p(\mathbf{v})$  as the M-step moves further and further away from the value  $\theta^{(0)}$  used in the E-step. Fortunately, the E-step reduces the gap to zero again as we enter the loop for the next time.

The EM algorithm contains a few different insights. First, there is the basic structure of the learning process, in which we update the model parameters to improve the likelihood of a completed dataset, where all missing variables have their values provided by an estimate of the posterior distribution. This particular insight is not unique to the EM algorithm. For example, using gradient descent to maximize the log-likelihood also has this same property; the log-likelihood gradient computations require taking expectations with respect to the posterior distribution over the hidden units. Another key insight in the EM algorithm is that we can continue to use one value of  $q$  even after we have moved to a different value of  $\theta$ . This particular insight is used throughout classical machine learning to derive large M-step updates. In the context of deep learning, most models are too complex to admit a tractable solution for an optimal large M-step update, so this second insight which is more unique to the EM algorithm is rarely used.

### 19.3 MAP Inference and Sparse Coding

We usually use the term inference to refer to computing the probability distribution over one set of variables given another. When training probabilistic models with latent variables, we are usually interested in computing  $p(\mathbf{h} | \mathbf{v})$ . An alternative form of inference is to compute the single most likely value of the missing variables, rather than to infer the entire distribution over their possible values. In the context

of latent variable models, this means computing

$$\mathbf{h}^* = \arg \max_{\mathbf{h}} p(\mathbf{h} | \mathbf{v}). \quad (19.9)$$

This is known as *maximum a posteriori* inference, abbreviated MAP inference.

MAP inference is usually not thought of as approximate inference—it does compute the exact most likely value of  $\mathbf{h}^*$ . However, if we wish to develop a learning process based on maximizing  $\mathcal{L}(\mathbf{v}, \mathbf{h}, q)$ , then it is helpful to think of MAP inference as a procedure that provides a value of  $q$ . In this sense, we can think of MAP inference as approximate inference, because it does not provide the optimal  $q$ .

Recall from Sec. 19.1 that exact inference consists of maximizing

$$\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q) = \mathbb{E}_{\mathbf{h} \sim q} [\log p(\mathbf{h}, \mathbf{v})] + H(q) \quad (19.10)$$

with respect to  $q$  over an unrestricted family of probability distributions, using an exact optimization algorithm. We can derive MAP inference as a form of approximate inference by restricting the family of distributions  $q$  may be drawn from. Specifically, we require  $q$  to take on a Dirac distribution:

$$q(\mathbf{h} | \mathbf{v}) = \delta(\mathbf{h} - \boldsymbol{\mu}). \quad (19.11)$$

This means that we can now control  $q$  entirely via  $\boldsymbol{\mu}$ . Dropping terms of  $\mathcal{L}$  that do not vary with  $\boldsymbol{\mu}$ , we are left with the optimization problem

$$\boldsymbol{\mu}^* = \max_{\boldsymbol{\mu}} \log p(\mathbf{h} = \boldsymbol{\mu}, \mathbf{v}), \quad (19.12)$$

which is equivalent to the MAP inference problem

$$\mathbf{h}^* = \max_{\mathbf{h}} p(\mathbf{h} | \mathbf{v}). \quad (19.13)$$

We can thus justify a learning procedure similar to EM, in which we alternate between performing MAP inference to infer  $\mathbf{h}^*$  and then update  $\boldsymbol{\theta}$  to increase  $\log p(\mathbf{h}^*, \mathbf{v})$ . As with EM, this is a form of coordinate ascent on  $\mathcal{L}$ , where we alternate between using inference to optimize  $\mathcal{L}$  with respect to  $q$  and using parameter updates to optimize  $\mathcal{L}$  with respect to  $\boldsymbol{\theta}$ . The procedure as a whole can be justified by the fact that  $\mathcal{L}$  is a lower bound on  $\log p(\mathbf{v})$ . In the case of MAP inference, this justification is rather vacuous, because the bound is infinitely loose, due to the Dirac distribution's differential entropy of negative infinity. However, adding noise to  $\boldsymbol{\mu}$  would make the bound meaningful again.

MAP inference is commonly used in deep learning as both a feature extractor and a learning mechanism. It is primarily used for sparse coding models.

Recall from Sec. 13.4 that sparse coding is a linear factor model that imposes a sparsity-inducing prior on its hidden units. A common choice is a factorial Laplace prior, with

$$p(h_i) = \frac{\lambda}{2} e^{-\frac{1}{2}\lambda|h_i|}. \quad (19.14)$$

The visible units are then generated by performing a linear transformation and adding noise:

$$p(\mathbf{x} | \mathbf{h}) = \mathcal{N}(\mathbf{x}; \mathbf{W}\mathbf{h} + \mathbf{b}, \beta\mathbf{I}). \quad (19.15)$$

Computing or even representing  $p(\mathbf{h} | \mathbf{v})$  is difficult. Every pair of variables  $h_i$  and  $h_j$  are both parents of  $\mathbf{v}$ . This means that when  $\mathbf{v}$  is observed, the graphical model contains an active path connecting  $h_i$  and  $h_j$ . All of the hidden units thus participate in one massive clique in  $p(\mathbf{h} | \mathbf{v})$ . If the model were Gaussian then these interactions could be modeled efficiently via the covariance matrix, but the sparse prior makes these interactions non-Gaussian.

Because  $p(\mathbf{h} | \mathbf{v})$  is intractable, so is the computation of the log-likelihood and its gradient. We thus cannot use exact maximum likelihood learning. Instead, we use MAP inference and learn the parameters by maximizing the ELBO defined by the Dirac distribution around the MAP estimate of  $\mathbf{h}$ .

If we concatenate all of the  $\mathbf{h}$  vectors in the training set into a matrix  $\mathbf{H}$ , then the sparse coding learning process consists of minimizing

$$J(\mathbf{H}, \mathbf{W}) = \sum_{i,j} |H_{i,j}| + \sum_{i,j} (\mathbf{X} - \mathbf{HW}^\top)_{i,j}^2. \quad (19.16)$$

Most applications of sparse coding also involve weight decay or a constraint on the norms of the columns of  $\mathbf{W}$ , in order to prevent the pathological solution with extremely small  $\mathbf{H}$  and large  $\mathbf{W}$ .

We can minimize  $J$  by alternating between minimization with respect to  $\mathbf{H}$  and minimization with respect to  $\mathbf{W}$ . Both sub-problems are convex. In fact, the minimization with respect to  $\mathbf{W}$  is just a linear regression problem. However, minimization of  $J$  with respect to both arguments is usually not a convex problem.

Minimization with respect to  $\mathbf{H}$  requires specialized algorithms such as the feature-sign search algorithm (Lee *et al.*, 2007).

## 19.4 Variational Inference and Learning

We have seen how the evidence lower bound  $\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q)$  is a lower bound on  $\log p(\mathbf{v}; \boldsymbol{\theta})$ , how inference can be viewed as maximizing  $\mathcal{L}$  with respect to  $q$ , and how learning can be viewed as maximizing  $\mathcal{L}$  with respect to  $\boldsymbol{\theta}$ . We have seen that the EM algorithm allows us to make large learning steps with a fixed  $q$  and that learning algorithms based on MAP inference allow us to learn using a point estimate of  $p(\mathbf{h} | \mathbf{v})$  rather than inferring the entire distribution. Now we develop the more general approach to variational learning.

The core idea behind variational learning is that we can maximize  $\mathcal{L}$  over a restricted family of distributions  $q$ . This family should be chosen so that it is easy to compute  $\mathbb{E}_q \log p(\mathbf{h}, \mathbf{v})$ . A typical way to do this is to introduce assumptions about how  $q$  factorizes.

A common approach to variational learning is to impose the restriction that  $q$  is a factorial distribution:

$$q(\mathbf{h} | \mathbf{v}) = \prod_i q(h_i | \mathbf{v}). \quad (19.17)$$

This is called the *mean field* approach. More generally, we can impose any graphical model structure we choose on  $q$ , to flexibly determine how many interactions we want our approximation to capture. This fully general graphical model approach is called *structured variational inference* (Saul and Jordan, 1996).

The beauty of the variational approach is that we do not need to specify a specific parametric form for  $q$ . We specify how it should factorize, but then the optimization problem determines the optimal probability distribution within those factorization constraints. For discrete latent variables, this just means that we use traditional optimization techniques to optimize a finite number of variables describing the  $q$  distribution. For continuous latent variables, this means that we use a branch of mathematics called calculus of variations to perform optimization over a space of functions, and actually determine which function should be used to represent  $q$ . Calculus of variations is the origin of the names “variational learning” and “variational inference,” though these names apply even when the latent variables are discrete and calculus of variations is not needed. In the case of continuous latent variables, calculus of variations is a powerful technique that removes much of the responsibility from the human designer of the model, who now must specify only how  $q$  factorizes, rather than needing to guess how to design a specific  $q$  that can accurately approximate the posterior.

Because  $\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q)$  is defined to be  $\log p(\mathbf{v}; \boldsymbol{\theta}) - D_{\text{KL}}(q(\mathbf{h} | \mathbf{v}) \| p(\mathbf{h} | \mathbf{v}; \boldsymbol{\theta}))$ , we can think of maximizing  $\mathcal{L}$  with respect to  $q$  as minimizing  $D_{\text{KL}}(q(\mathbf{h} | \mathbf{v}) \| p(\mathbf{h} | \mathbf{v}))$ .

In this sense, we are fitting  $q$  to  $p$ . However, we are doing so with the opposite direction of the KL divergence than we are used to using for fitting an approximation. When we use maximum likelihood learning to fit a model to data, we minimize  $D_{\text{KL}}(p_{\text{data}} \| p_{\text{model}})$ . As illustrated in Fig. 3.6, this means that maximum likelihood encourages the model to have high probability everywhere that the data has high probability, while our optimization-based inference procedure encourages  $q$  to have low probability everywhere the true posterior has low probability. Both directions of the KL divergence can have desirable and undesirable properties. The choice of which to use depends on which properties are the highest priority for each application. In the case of the inference optimization problem, we choose to use  $D_{\text{KL}}(q(\mathbf{h} | \mathbf{v}) \| p(\mathbf{h} | \mathbf{v}))$  for computational reasons. Specifically, computing  $D_{\text{KL}}(q(\mathbf{h} | \mathbf{v}) \| p(\mathbf{h} | \mathbf{v}))$  involves evaluating expectations with respect to  $q$ , so by designing  $q$  to be simple, we can simplify the required expectations. The opposite direction of the KL divergence would require computing expectations with respect to the true posterior. Because the form of the true posterior is determined by the choice of model, we cannot design a reduced-cost approach to computing  $D_{\text{KL}}(p(\mathbf{h} | \mathbf{v}) \| q(\mathbf{h} | \mathbf{v}))$  exactly.

### 19.4.1 Discrete Latent Variables

Variational inference with discrete latent variables is relatively straightforward. We define a distribution  $q$ , typically one where each factor of  $q$  is just defined by a lookup table over discrete states. In the simplest case,  $\mathbf{h}$  is binary and we make the mean field assumption that  $q$  factorizes over each individual  $h_i$ . In this case we can parametrize  $q$  with a vector  $\hat{\mathbf{h}}$  whose entries are probabilities. Then  $q(h_i = 1 | \mathbf{v}) = \hat{h}_i$ .

After determining how to represent  $q$ , we simply optimize its parameters. In the case of discrete latent variables, this is just a standard optimization problem. In principle the selection of  $q$  could be done with any optimization algorithm, such as gradient descent.

Because this optimization must occur in the inner loop of a learning algorithm, it must be very fast. To achieve this speed, we typically use special optimization algorithms that are designed to solve comparatively small and simple problems in very few iterations. A popular choice is to iterate fixed point equations, in other words, to solve

$$\frac{\partial}{\partial \hat{h}_i} \mathcal{L} = 0 \quad (19.18)$$

for  $\hat{h}_i$ . We repeatedly update different elements of  $\hat{\mathbf{h}}$  until we satisfy a convergence

criterion.

To make this more concrete, we show how to apply variational inference to the *binary sparse coding* model (we present here the model developed by Henniges *et al.* (2010) but demonstrate traditional, generic mean field applied to the model, while they introduce a specialized algorithm). This derivation goes into considerable mathematical detail and is intended for the reader who wishes to fully resolve any ambiguity in the high-level conceptual description of variational inference and learning we have presented so far. Readers who do not plan to derive or implement variational learning algorithms may safely skip to the next section without missing any new high-level concepts. Readers who proceed with the binary sparse coding example are encouraged to review the list of useful properties of functions that commonly arise in probabilistic models in Sec. 3.10. We use these properties liberally throughout the following derivations without highlighting exactly where we use each one.

In the binary sparse coding model, the input  $\mathbf{v} \in \mathbb{R}^n$  is generated from the model by adding Gaussian noise to the sum of  $m$  different components which can each be present or absent. Each component is switched on or off by the corresponding hidden unit in  $\mathbf{h} \in \{0, 1\}^m$ :

$$p(h_i = 1) = \sigma(b_i) \quad (19.19)$$

$$p(\mathbf{v} | \mathbf{h}) = \mathcal{N}(\mathbf{v}; \mathbf{W}\mathbf{h}, \boldsymbol{\beta}^{-1}) \quad (19.20)$$

where  $\mathbf{b}$  is a learnable set of biases,  $\mathbf{W}$  is a learnable weight matrix, and  $\boldsymbol{\beta}$  is a learnable, diagonal precision matrix.

Training this model with maximum likelihood requires taking the derivative with respect to the parameters. Consider the derivative with respect to one of the biases:

$$\frac{\partial}{\partial b_i} \log p(\mathbf{v}) \quad (19.21)$$

$$= \frac{\frac{\partial}{\partial b_i} p(\mathbf{v})}{p(\mathbf{v})} \quad (19.22)$$

$$= \frac{\frac{\partial}{\partial b_i} \sum_{\mathbf{h}} p(\mathbf{h}, \mathbf{v})}{p(\mathbf{v})} \quad (19.23)$$

$$= \frac{\frac{\partial}{\partial b_i} \sum_{\mathbf{h}} p(\mathbf{h}) p(\mathbf{v} | \mathbf{h})}{p(\mathbf{v})} \quad (19.24)$$

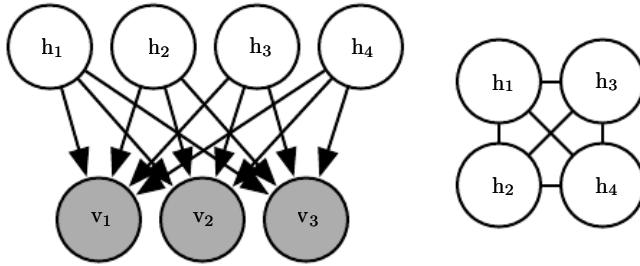


Figure 19.2: The graph structure of a binary sparse coding model with four hidden units. (Left) The graph structure of  $p(\mathbf{h}, \mathbf{v})$ . Note that the edges are directed, and that every two hidden units are co-parents of every visible unit. (Right) The graph structure of  $p(\mathbf{h} | \mathbf{v})$ . In order to account for the active paths between co-parents, the posterior distribution needs an edge between all of the hidden units.

$$= \frac{\sum_{\mathbf{h}} p(\mathbf{v} | \mathbf{h}) \frac{\partial}{\partial b_i} p(\mathbf{h})}{p(\mathbf{v})} \quad (19.25)$$

$$= \sum_{\mathbf{h}} p(\mathbf{h} | \mathbf{v}) \frac{\frac{\partial}{\partial b_i} p(\mathbf{h})}{p(\mathbf{h})} \quad (19.26)$$

$$= \mathbb{E}_{\mathbf{h} \sim p(\mathbf{h} | \mathbf{v})} \frac{\partial}{\partial b_i} \log p(\mathbf{h}). \quad (19.27)$$

This requires computing expectations with respect to  $p(\mathbf{h} | \mathbf{v})$ . Unfortunately,  $p(\mathbf{h} | \mathbf{v})$  is a complicated distribution. See Fig. 19.2 for the graph structure of  $p(\mathbf{h}, \mathbf{v})$  and  $p(\mathbf{h} | \mathbf{v})$ . The posterior distribution corresponds to the complete graph over the hidden units, so variable elimination algorithms do not help us to compute the required expectations any faster than brute force.

We can resolve this difficulty by using variational inference and variational learning instead.

We can make a mean field approximation:

$$q(\mathbf{h} | \mathbf{v}) = \prod_i q(h_i | \mathbf{v}). \quad (19.28)$$

The latent variables of the binary sparse coding model are binary, so to represent a factorial  $q$  we simply need to model  $m$  Bernoulli distributions  $q(h_i | \mathbf{v})$ . A natural way to represent the means of the Bernoulli distributions is with a vector  $\hat{\mathbf{h}}$  of probabilities, with  $q(h_i = 1 | \mathbf{v}) = \hat{h}_i$ . We impose a restriction that  $\hat{h}_i$  is never equal to 0 or to 1, in order to avoid errors when computing, for example,  $\log \hat{h}_i$ .

We will see that the variational inference equations never assign 0 or 1 to  $\hat{h}_i$

analytically. However, in a software implementation, machine rounding error could result in 0 or 1 values. In software, we may wish to implement binary sparse coding using an unrestricted vector of variational parameters  $\mathbf{z}$  and obtain  $\hat{\mathbf{h}}$  via the relation  $\hat{\mathbf{h}} = \sigma(\mathbf{z})$ . We can thus safely compute  $\log \hat{h}_i$  on a computer by using the identity  $\log \sigma(z_i) = -\zeta(-z_i)$  relating the sigmoid and the softplus.

To begin our derivation of variational learning in the binary sparse coding model, we show that the use of this mean field approximation makes learning tractable.

The evidence lower bound is given by

$$\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q) \tag{19.29}$$

$$= \mathbb{E}_{\mathbf{h} \sim q} [\log p(\mathbf{h}, \mathbf{v})] + H(q) \tag{19.30}$$

$$= \mathbb{E}_{\mathbf{h} \sim q} [\log p(\mathbf{h}) + \log p(\mathbf{v} \mid \mathbf{h}) - \log q(\mathbf{h} \mid \mathbf{v})] \tag{19.31}$$

$$= \mathbb{E}_{\mathbf{h} \sim q} \left[ \sum_{i=1}^m \log p(h_i) + \sum_{i=1}^n \log p(v_i \mid \mathbf{h}) - \sum_{i=1}^m \log q(h_i \mid \mathbf{v}) \right] \tag{19.32}$$

$$= \sum_{i=1}^m \left[ \hat{h}_i (\log \sigma(b_i) - \log \hat{h}_i) + (1 - \hat{h}_i) (\log \sigma(-b_i) - \log(1 - \hat{h}_i)) \right] \tag{19.33}$$

$$+ \mathbb{E}_{\mathbf{h} \sim q} \left[ \sum_{i=1}^n \log \sqrt{\frac{\beta_i}{2\pi}} \exp \left( -\frac{\beta_i}{2} (v_i - \mathbf{W}_{i,:} \mathbf{h})^2 \right) \right] \tag{19.34}$$

$$= \sum_{i=1}^m \left[ \hat{h}_i (\log \sigma(b_i) - \log \hat{h}_i) + (1 - \hat{h}_i) (\log \sigma(-b_i) - \log(1 - \hat{h}_i)) \right] \tag{19.35}$$

$$+ \frac{1}{2} \sum_{i=1}^n \left[ \log \frac{\beta_i}{2\pi} - \beta_i \left( v_i^2 - 2v_i \mathbf{W}_{i,:} \hat{\mathbf{h}} + \sum_j \left[ W_{i,j}^2 \hat{h}_j + \sum_{k \neq j} W_{i,j} W_{i,k} \hat{h}_j \hat{h}_k \right] \right) \right]. \tag{19.36}$$

While these equations are somewhat unappealing aesthetically, they show that  $\mathcal{L}$  can be expressed in a small number of simple arithmetic operations. The evidence lower bound  $\mathcal{L}$  is therefore tractable. We can use  $\mathcal{L}$  as a replacement for the intractable log-likelihood.

In principle, we could simply run gradient ascent on both  $\mathbf{v}$  and  $\mathbf{h}$  and this would make a perfectly acceptable combined inference and training algorithm. Usually, however, we do not do this, for two reasons. First, this would require storing  $\hat{\mathbf{h}}$  for each  $\mathbf{v}$ . We typically prefer algorithms that do not require per-example memory. It is difficult to scale learning algorithms to billions of examples if we must remember a dynamically updated vector associated with each example.

Second, we would like to be able to extract the features  $\hat{\mathbf{h}}$  very quickly, in order to recognize the content of  $\mathbf{v}$ . In a realistic deployed setting, we would need to be able to compute  $\hat{\mathbf{h}}$  in real time.

For both these reasons, we typically do not use gradient descent to compute the mean field parameters  $\hat{\mathbf{h}}$ . Instead, we rapidly estimate them with fixed point equations.

The idea behind fixed point equations is that we are seeking a local maximum with respect to  $\hat{\mathbf{h}}$ , where  $\nabla_{\hat{\mathbf{h}}} \mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, \hat{\mathbf{h}}) = \mathbf{0}$ . We cannot efficiently solve this equation with respect to all of  $\hat{\mathbf{h}}$  simultaneously. However, we can solve for a single variable:

$$\frac{\partial}{\partial \hat{h}_i} \mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, \hat{\mathbf{h}}) = 0. \quad (19.37)$$

We can then iteratively apply the solution to the equation for  $i = 1, \dots, m$ , and repeat the cycle until we satisfy a converge criterion. Common convergence criteria include stopping when a full cycle of updates does not improve  $\mathcal{L}$  by more than some tolerance amount, or when the cycle does not change  $\hat{\mathbf{h}}$  by more than some amount.

Iterating mean field fixed point equations is a general technique that can provide fast variational inference in a broad variety of models. To make this more concrete, we show how to derive the updates for the binary sparse coding model in particular.

First, we must write an expression for the derivatives with respect to  $\hat{h}_i$ . To do so, we substitute Eq. 19.36 into the left side of Eq. 19.37:

$$\frac{\partial}{\partial \hat{h}_i} \mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, \hat{\mathbf{h}}) \quad (19.38)$$

$$= \frac{\partial}{\partial \hat{h}_i} \left[ \sum_{j=1}^m \left[ \hat{h}_j (\log \sigma(b_j) - \log \hat{h}_j) + (1 - \hat{h}_j) (\log \sigma(-b_j) - \log(1 - \hat{h}_j)) \right] \right] \quad (19.39)$$

$$+ \frac{1}{2} \sum_{j=1}^n \left[ \log \frac{\beta_j}{2\pi} - \beta_j \left( v_j^2 - 2v_j \mathbf{W}_{j,:} \hat{\mathbf{h}} - \sum_k \left[ W_{j,k}^2 \hat{h}_k + \sum_{l \neq k} W_{j,k} W_{j,l} \hat{h}_k \hat{h}_l \right] \right) \right] \quad (19.40)$$

$$= \log \sigma(b_i) - \log \hat{h}_i - 1 + \log(1 - \hat{h}_i) + 1 - \log \sigma(-b_i) \quad (19.41)$$

$$+ \sum_{j=1}^n \left[ \beta_j \left( v_j W_{j,i} - \frac{1}{2} W_{j,i}^2 - \sum_{k \neq i} W_{j,k} W_{j,i} \hat{h}_k \right) \right] \quad (19.42)$$

$$= b_i - \log \hat{h}_i + \log(1 - \hat{h}_i) + \mathbf{v}^\top \boldsymbol{\beta} \mathbf{W}_{:,i} - \frac{1}{2} \mathbf{W}_{:,i} \boldsymbol{\beta} \mathbf{W}_{:,i} - \sum_{j \neq i} \mathbf{W}_{:,j}^\top \boldsymbol{\beta} \mathbf{W}_{:,i} \hat{h}_j. \quad (19.43)$$

To apply the fixed point update inference rule, we solve for the  $\hat{h}_i$  that sets Eq. 19.43 to 0:

$$\hat{h}_i = \sigma \left( b_i + \mathbf{v}^\top \boldsymbol{\beta} \mathbf{W}_{:,i} - \frac{1}{2} \mathbf{W}_{:,i} \boldsymbol{\beta} \mathbf{W}_{:,i} - \sum_{j \neq i} \mathbf{W}_{:,j}^\top \boldsymbol{\beta} \mathbf{W}_{:,i} \hat{h}_j \right). \quad (19.44)$$

At this point, we can see that there is a close connection between recurrent neural networks and inference in graphical models. Specifically, the mean field fixed point equations defined a recurrent neural network. The task of this network is to perform inference. We have described how to derive this network from a model description, but it is also possible to train the inference network directly. Several ideas based on this theme are described in Chapter 20.

In the case of binary sparse coding, we can see that the recurrent network connection specified by Eq. 19.44 consists of repeatedly updating the hidden units based on the changing values of the neighboring hidden units. The input always sends a fixed message of  $\mathbf{v}^\top \boldsymbol{\beta} \mathbf{W}$  to the hidden units, but the hidden units constantly update the message they send to each other. Specifically, two units  $\hat{h}_i$  and  $\hat{h}_j$  inhibit each other when their weight vectors are aligned. This is a form of competition—between two hidden units that both explain the input, only the one that explains the input best will be allowed to remain active. This competition is the mean field approximation’s attempt to capture the explaining away interactions in the binary sparse coding posterior. The explaining away effect actually should cause a multi-modal posterior, so that if we draw samples from the posterior, some samples will have one unit active, other samples will have the other unit active, but very few samples have both active. Unfortunately, explaining away interactions cannot be modeled by the factorial  $q$  used for mean field, so the mean field approximation is forced to choose one mode to model. This is an instance of the behavior illustrated in Fig. 3.6.

We can rewrite Eq. 19.44 into an equivalent form that reveals some further insights:

$$\hat{h}_i = \sigma \left( b_i + \left( \mathbf{v} - \sum_{j \neq i} \mathbf{W}_{:,j} \hat{h}_j \right)^\top \boldsymbol{\beta} \mathbf{W}_{:,i} - \frac{1}{2} \mathbf{W}_{:,i} \boldsymbol{\beta} \mathbf{W}_{:,i} \right). \quad (19.45)$$

In this reformulation, we see the input at each step as consisting of  $\mathbf{v} - \sum_{j \neq i} \mathbf{W}_{:,j} \hat{h}_j$  rather than  $\mathbf{v}$ . We can thus think of unit  $i$  as attempting to encode the residual

error in  $\mathbf{v}$  given the code of the other units. We can thus think of sparse coding as an iterative autoencoder, that repeatedly encodes and decodes its input, attempting to fix mistakes in the reconstruction after each iteration.

In this example, we have derived an update rule that updates a single unit at a time. It would be advantageous to be able to update more units simultaneously. Some graphical models, such as deep Boltzmann machines, are structured in such a way that we can solve for many entries of  $\hat{\mathbf{h}}$  simultaneously. Unfortunately, binary sparse coding does not admit such block updates. Instead, we can use a heuristic technique called *damping* to perform block updates. In the damping approach, we solve for the individually optimal values of every element of  $\hat{\mathbf{h}}$ , then move all of the values in a small step in that direction. This approach is no longer guaranteed to increase  $\mathcal{L}$  at each step, but works well in practice for many models. See [Koller and Friedman \(2009\)](#) for more information about choosing the degree of synchrony and damping strategies in message passing algorithms.

#### 19.4.2 Calculus of Variations

Before continuing with our presentation of variational learning, we must briefly introduce an important set of mathematical tools used in variational learning: *calculus of variations*.

Many machine learning techniques are based on minimizing a function  $J(\boldsymbol{\theta})$  by finding the input vector  $\boldsymbol{\theta} \in \mathbb{R}^n$  for which it takes on its minimal value. This can be accomplished with multivariate calculus and linear algebra, by solving for the critical points where  $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbf{0}$ . In some cases, we actually want to solve for a function  $f(\mathbf{x})$ , such as when we want to find the probability density function over some random variable. This is what calculus of variations enables us to do.

A function of a function  $f$  is known as a *functional*  $J[f]$ . Much as we can take partial derivatives of a function with respect to elements of its vector-valued argument, we can take *functional derivatives*, also known as *variational derivatives*, of a functional  $J[f]$  with respect to individual values of the function  $f(\mathbf{x})$  at any specific value of  $\mathbf{x}$ . The functional derivative of the functional  $J$  with respect to the value of the function  $f$  at point  $\mathbf{x}$  is denoted  $\frac{\delta}{\delta f(\mathbf{x})} J$ .

A complete formal development of functional derivatives is beyond the scope of this book. For our purposes, it is sufficient to state that for differentiable functions  $f(\mathbf{x})$  and differentiable functions  $g(y, \mathbf{x})$  with continuous derivatives, that

$$\frac{\delta}{\delta f(\mathbf{x})} \int g(f(\mathbf{x}), \mathbf{x}) d\mathbf{x} = \frac{\partial}{\partial y} g(f(\mathbf{x}), \mathbf{x}). \quad (19.46)$$

To gain some intuition for this identity, one can think of  $f(\mathbf{x})$  as being a vector with uncountably many elements, indexed by a real vector  $\mathbf{x}$ . In this (somewhat incomplete view), the identity providing the functional derivatives is the same as we would obtain for a vector  $\boldsymbol{\theta} \in \mathbb{R}^n$  indexed by positive integers:

$$\frac{\partial}{\partial \theta_i} \sum_j g(\theta_j, j) = \frac{\partial}{\partial \theta_i} g(\theta_i, i). \quad (19.47)$$

Many results in other machine learning publications are presented using the more general *Euler-Lagrange equation* which allows  $g$  to depend on the derivatives of  $f$  as well as the value of  $f$ , but we do not need this fully general form for the results presented in this book.

To optimize a function with respect to a vector, we take the gradient of the function with respect to the vector and solve for the point where every element of the gradient is equal to zero. Likewise, we can optimize a functional by solving for the function where the functional derivative at every point is equal to zero.

As an example of how this process works, consider the problem of finding the probability distribution function over  $x \in \mathbb{R}$  that has maximal differential entropy. Recall that the entropy of a probability distribution  $p(x)$  is defined as

$$H[p] = -\mathbb{E}_x \log p(x). \quad (19.48)$$

For continuous values, the expectation is an integral:

$$H[p] = - \int p(x) \log p(x) dx. \quad (19.49)$$

We cannot simply maximize  $H(x)$  with respect to the function  $p(x)$ , because the result might not be a probability distribution. Instead, we need to use Lagrange multipliers, to add a constraint that  $p(x)$  integrates to 1. Also, the entropy increases without bound as the variance increases. This makes the question of which distribution has the greatest entropy uninteresting. Instead, we ask which distribution has maximal entropy for fixed variance  $\sigma^2$ . Finally, the problem is underdetermined because the distribution can be shifted arbitrarily without changing the entropy. To impose a unique solution, we add a constraint that the mean of the distribution be  $\mu$ . The Lagrangian functional for this optimization problem is

$$\mathcal{L}[p] = \lambda_1 \left( \int p(x) dx - 1 \right) + \lambda_2 (\mathbb{E}[x] - \mu) + \lambda_3 (\mathbb{E}[(x - \mu)^2] - \sigma^2) + H[p] \quad (19.50)$$

$$= \int (\lambda_1 p(x) + \lambda_2 p(x)x + \lambda_3 p(x)(x - \mu)^2 - p(x) \log p(x)) dx - \lambda_1 - \mu\lambda_2 - \sigma^2 \lambda_3. \quad (19.51)$$

To minimize the Lagrangian with respect to  $p$ , we set the functional derivatives equal to 0:

$$\forall x, \frac{\delta}{\delta p(x)} \mathcal{L} = \lambda_1 + \lambda_2 x + \lambda_3(x - \mu)^2 - 1 - \log p(x) = 0. \quad (19.52)$$

This condition now tells us the functional form of  $p(x)$ . By algebraically re-arranging the equation, we obtain

$$p(x) = \exp(-\lambda_1 - \lambda_2 x + \lambda_3(x - \mu)^2 + 1). \quad (19.53)$$

We never assumed directly that  $p(x)$  would take this functional form; we obtained the expression itself by analytically minimizing a functional. To finish the minimization problem, we must choose the  $\lambda$  values to ensure that all of our constraints are satisfied. We are free to choose any  $\lambda$  values, because the gradient of the Lagrangian with respect to the  $\lambda$  variables is zero so long as the constraints are satisfied. To satisfy all of the constraints, we may set  $\lambda_1 = \log \sigma \sqrt{2\pi}$ ,  $\lambda_2 = 0$ , and  $\lambda_3 = -\frac{1}{2\sigma^2}$  to obtain

$$p(x) = \mathcal{N}(x; \mu, \sigma^2). \quad (19.54)$$

This is one reason for using the normal distribution when we do not know the true distribution. Because the normal distribution has the maximum entropy, we impose the least possible amount of structure by making this assumption.

While examining the critical points of the Lagrangian functional for the entropy, we found only one critical point, corresponding to maximizing the entropy for fixed variance. What about the probability distribution function that **minimizes** the entropy? Why did we not find a second critical point corresponding to the minimum? The reason is that there is no specific function that achieves minimal entropy. As functions place more probability density on the two points  $x = \mu + \sigma$  and  $x = \mu - \sigma$ , and place less probability density on all other values of  $x$ , they lose entropy while maintaining the desired variance. However, any function placing exactly zero mass on all but two points does not integrate to one, and is not a valid probability distribution. There thus is no single minimal entropy probability distribution function, much as there is no single minimal positive real number. Instead, we can say that there is a sequence of probability distributions converging toward putting mass only on these two points. This degenerate scenario may be described as a mixture of Dirac distributions. Because Dirac distributions are not described by a single probability distribution function, no Dirac or mixture of

Dirac distribution corresponds to a single specific point in function space. These distributions are thus invisible to our method of solving for a specific point where the functional derivatives are zero. This is a limitation of the method. Distributions such as the Dirac must be found by other methods, such as guessing the solution and then proving that it is correct.

### 19.4.3 Continuous Latent Variables

When our graphical model contains continuous latent variables, we may still perform variational inference and learning by maximizing  $\mathcal{L}$ . However, we must now use calculus of variations when maximizing  $\mathcal{L}$  with respect to  $q(\mathbf{h} | \mathbf{v})$ .

In most cases, practitioners need not solve any calculus of variations problems themselves. Instead, there is a general equation for the mean field fixed point updates. If we make the mean field approximation

$$q(\mathbf{h} | \mathbf{v}) = \prod_i q(h_i | \mathbf{v}), \quad (19.55)$$

and fix  $q(h_j | \mathbf{v})$  for all  $j \neq i$ , then the optimal  $q(h_i | \mathbf{v})$  may be obtained by normalizing the unnormalized distribution

$$\tilde{q}(h_i | \mathbf{v}) = \exp (\mathbb{E}_{\mathbf{h}_{-i} \sim q(\mathbf{h}_{-i} | \mathbf{v})} \log \tilde{p}(\mathbf{v}, \mathbf{h})) \quad (19.56)$$

so long as  $p$  does not assign 0 probability to any joint configuration of variables. Carrying out the expectation inside the equation will yield the correct functional form of  $q(h_i | \mathbf{v})$ . It is only necessary to derive functional forms of  $q$  directly using calculus of variations if one wishes to develop a new form of variational learning; Eq. 19.56 yields the mean field approximation for any probabilistic model.

Eq. 19.56 is a fixed point equation, designed to be iteratively applied for each value of  $i$  repeatedly until convergence. However, it also tells us more than that. It tells us the functional form that the optimal solution will take, whether we arrive there by fixed point equations or not. This means we can take the functional form from that equation but regard some of the values that appear in it as parameters, that we can optimize with any optimization algorithm we like.

As an example, consider a very simple probabilistic model, with latent variables  $\mathbf{h} \in \mathbb{R}^2$  and just one visible variable,  $v$ . Suppose that  $p(\mathbf{h}) = \mathcal{N}(\mathbf{h}; \mathbf{0}, \mathbf{I})$  and  $p(v | \mathbf{h}) = \mathcal{N}(v; \mathbf{w}^\top \mathbf{h}; 1)$ . We could actually simplify this model by integrating out  $\mathbf{h}$ ; the result is just a Gaussian distribution over  $v$ . The model itself is not interesting; we have constructed it only to provide a simple demonstration of how calculus of variations may be applied to probabilistic modeling.

The true posterior is given, up to a normalizing constant, by

$$p(\mathbf{h} | \mathbf{v}) \quad (19.57)$$

$$\propto p(\mathbf{h}, \mathbf{v}) \quad (19.58)$$

$$= p(h_1)p(h_2)p(\mathbf{v} | \mathbf{h}) \quad (19.59)$$

$$\propto \exp\left(-\frac{1}{2} [h_1^2 + h_2^2 + (v - h_1 w_1 - h_2 w_2)^2]\right) \quad (19.60)$$

$$= \exp\left(-\frac{1}{2} [h_1^2 + h_2^2 + v^2 + h_1^2 w_1^2 + h_2^2 w_2^2 - 2vh_1 w_1 - 2vh_2 w_2 - 2h_1 w_1 h_2 w_2]\right). \quad (19.61)$$

Due to the presence of the terms multiplying  $h_1$  and  $h_2$  together, we can see that the true posterior does not factorize over  $h_1$  and  $h_2$ .

Applying Eq. 19.56, we find that

$$\tilde{q}(h_1 | \mathbf{v}) \quad (19.62)$$

$$= \exp(\mathbb{E}_{h_2 \sim q(h_2 | \mathbf{v})} \log \tilde{p}(\mathbf{v}, \mathbf{h})) \quad (19.63)$$

$$= \exp\left(-\frac{1}{2} \mathbb{E}_{h_2 \sim q(h_2 | \mathbf{v})} [h_1^2 + h_2^2 + v^2 + h_1^2 w_1^2 + h_2^2 w_2^2 \quad (19.64)\right.$$

$$\left. - 2vh_1 w_1 - 2vh_2 w_2 - 2h_1 w_1 h_2 w_2]\right). \quad (19.65)$$

From this, we can see that there are effectively only two values we need to obtain from  $q(h_2 | \mathbf{v})$ :  $\mathbb{E}_{h_2 \sim q(h_2 | \mathbf{v})}[h_2]$  and  $\mathbb{E}_{h_2 \sim q(h_2 | \mathbf{v})}[h_2^2]$ . Writing these as  $\langle h_2 \rangle$  and  $\langle h_2^2 \rangle$ , we obtain

$$\tilde{q}(h_1 | \mathbf{v}) = \exp\left(-\frac{1}{2} [h_1^2 + \langle h_2^2 \rangle + v^2 + h_1^2 w_1^2 + \langle h_2^2 \rangle w_2^2 \quad (19.66)\right.$$

$$\left. - 2vh_1 w_1 - 2v\langle h_2 \rangle w_2 - 2h_1 w_1 \langle h_2 \rangle w_2]\right). \quad (19.67)$$

From this, we can see that  $\tilde{q}$  has the functional form of a Gaussian. We can thus conclude  $q(\mathbf{h} | \mathbf{v}) = \mathcal{N}(\mathbf{h}; \boldsymbol{\mu}, \boldsymbol{\beta}^{-1})$  where  $\boldsymbol{\mu}$  and diagonal  $\boldsymbol{\beta}$  are variational parameters that we can optimize using any technique we choose. It is important to recall that we did not ever assume that  $q$  would be Gaussian; its Gaussian form was derived automatically by using calculus of variations to maximize  $q$  with respect to  $\mathcal{L}$ . Using the same approach on a different model could yield a different functional form of  $q$ .

This was of course, just a small case constructed for demonstration purposes. For examples of real applications of variational learning with continuous variables in the context of deep learning, see Goodfellow *et al.* (2013d).

### 19.4.4 Interactions between Learning and Inference

Using approximate inference as part of a learning algorithm affects the learning process, and this in turn affects the accuracy of the inference algorithm.

Specifically, the training algorithm tends to adapt the model in a way that makes the approximating assumptions underlying the approximate inference algorithm become more true. When training the parameters, variational learning increases

$$\mathbb{E}_{\mathbf{h} \sim q} \log p(\mathbf{v}, \mathbf{h}). \quad (19.68)$$

For a specific  $\mathbf{v}$ , this increases  $p(\mathbf{h} | \mathbf{v})$  for values of  $\mathbf{h}$  that have high probability under  $q(\mathbf{h} | \mathbf{v})$  and decreases  $p(\mathbf{h} | \mathbf{v})$  for values of  $\mathbf{h}$  that have low probability under  $q(\mathbf{h} | \mathbf{v})$ .

This behavior causes our approximating assumptions to become self-fulfilling prophecies. If we train the model with a unimodal approximate posterior, we will obtain a model with a true posterior that is far closer to unimodal than we would have obtained by training the model with exact inference.

Computing the true amount of harm imposed on a model by a variational approximation is thus very difficult. There exist several methods for estimating  $\log p(\mathbf{v})$ . We often estimate  $\log p(\mathbf{v}; \boldsymbol{\theta})$  after training the model, and find that the gap with  $\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q)$  is small. From this, we can conclude that our variational approximation is accurate for the specific value of  $\boldsymbol{\theta}$  that we obtained from the learning process. We should not conclude that our variational approximation is accurate in general or that the variational approximation did little harm to the learning process. To measure the true amount of harm induced by the variational approximation, we would need to know  $\boldsymbol{\theta}^* = \max_{\boldsymbol{\theta}} \log p(\mathbf{v}; \boldsymbol{\theta})$ . It is possible for  $\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q) \approx \log p(\mathbf{v}; \boldsymbol{\theta})$  and  $\log p(\mathbf{v}; \boldsymbol{\theta}) \ll \log p(\mathbf{v}; \boldsymbol{\theta}^*)$  to hold simultaneously. If  $\max_q \mathcal{L}(\mathbf{v}, \boldsymbol{\theta}^*, q) \ll \log p(\mathbf{v}; \boldsymbol{\theta}^*)$ , because  $\boldsymbol{\theta}^*$  induces too complicated of a posterior distribution for our  $q$  family to capture, then the learning process will never approach  $\boldsymbol{\theta}^*$ . Such a problem is very difficult to detect, because we can only know for sure that it happened if we have a superior learning algorithm that can find  $\boldsymbol{\theta}^*$  for comparison.

## 19.5 Learned Approximate Inference

We have seen that inference can be thought of as an optimization procedure that increases the value of a function  $\mathcal{L}$ . Explicitly performing optimization via iterative procedures such as fixed point equations or gradient-based optimization is often very expensive and time-consuming. Many approaches to inference avoid

this expense by learning to perform approximate inference. Specifically, we can think of the optimization process as a function  $f$  that maps an input  $\mathbf{v}$  to an approximate distribution  $q^* = \arg \max_q \mathcal{L}(\mathbf{v}, q)$ . Once we think of the multi-step iterative optimization process as just being a function, we can approximate it with a neural network that implements an approximation  $\hat{f}(\mathbf{v}; \boldsymbol{\theta})$ .

### 19.5.1 Wake-Sleep

One of the main difficulties with training a model to infer  $\mathbf{h}$  from  $\mathbf{v}$  is that we do not have a supervised training set with which to train the model. Given a  $\mathbf{v}$ , we do not know the appropriate  $\mathbf{h}$ . The mapping from  $\mathbf{v}$  to  $\mathbf{h}$  depends on the choice of model family, and evolves throughout the learning process as  $\boldsymbol{\theta}$  changes. The wake-sleep algorithm (Hinton *et al.*, 1995b; Frey *et al.*, 1996) resolves this problem by drawing samples of both  $\mathbf{h}$  and  $\mathbf{v}$  from the model distribution. For example, in a directed model, this can be done cheaply by performing ancestral sampling beginning at  $\mathbf{h}$  and ending at  $\mathbf{v}$ . The inference network can then be trained to perform the reverse mapping: predicting which  $\mathbf{h}$  caused the present  $\mathbf{v}$ . The main drawback to this approach is that we will only be able to train the inference network on values of  $\mathbf{v}$  that have high probability under the model. Early in learning, the model distribution will not resemble the data distribution, so the inference network will not have an opportunity to learn on samples that resemble data.

In Sec. 18.2 we saw that one possible explanation for the role of dream sleep in human beings and animals is that dreams could provide the negative phase samples that Monte Carlo training algorithms use to approximate the negative gradient of the log partition function of undirected models. Another possible explanation for biological dreaming is that it is providing samples from  $p(\mathbf{h}, \mathbf{v})$  which can be used to train an inference network to predict  $\mathbf{h}$  given  $\mathbf{v}$ . In some senses, this explanation is more satisfying than the partition function explanation. Monte Carlo algorithms generally do not perform well if they are run using only the positive phase of the gradient for several steps then with only the negative phase of the gradient for several steps. Human beings and animals are usually awake for several consecutive hours then asleep for several consecutive hours. It is not readily apparent how this schedule could support Monte Carlo training of an undirected model. Learning algorithms based on maximizing  $\mathcal{L}$  can be run with prolonged periods of improving  $q$  and prolonged periods of improving  $\boldsymbol{\theta}$ , however. If the role of biological dreaming is to train networks for predicting  $q$ , then this explains how animals are able to remain awake for several hours (the longer they are awake, the greater the gap between  $\mathcal{L}$  and  $\log p(\mathbf{v})$ , but  $\mathcal{L}$  will remain a lower bound) and to remain asleep

for several hours (the generative model itself is not modified during sleep) without damaging their internal models. Of course, these ideas are purely speculative, and there is no hard evidence to suggest that dreaming accomplishes either of these goals. Dreaming may also serve reinforcement learning rather than probabilistic modeling, by sampling synthetic experiences from the animal’s transition model, on which to train the animal’s policy. Or sleep may serve some other purpose not yet anticipated by the machine learning community.

### 19.5.2 Other Forms of Learned Inference

This strategy of learned approximate inference has also been applied to other models. [Salakhutdinov and Larochelle \(2010\)](#) showed that a single pass in a learned inference network could yield faster inference than iterating the mean field fixed point equations in a DBM. The training procedure is based on running the inference network, then applying one step of mean field to improve its estimates, and training the inference network to output this refined estimate instead of its original estimate.

We have already seen in Sec. 14.8 that the predictive sparse decomposition model trains a shallow encoder network to predict a sparse code for the input. This can be seen as a hybrid between an autoencoder and sparse coding. It is possible to devise probabilistic semantics for the model, under which the encoder may be viewed as performing learned approximate MAP inference. Due to its shallow encoder, PSD is not able to implement the kind of competition between units that we have seen in mean field inference. However, that problem can be remedied by training a deep encoder to perform learned approximate inference, as in the ISTA technique ([Gregor and LeCun, 2010b](#)).

Learned approximate inference has recently become one of the dominant approaches to generative modeling, in the form of the variational autoencoder ([Kingma, 2013; Rezende et al., 2014](#)). In this elegant approach, there is no need to construct explicit targets for the inference network. Instead, the inference network is simply used to define  $\mathcal{L}$ elegant approach, there is no need the inference network are adapted to increase  $\mathcal{L}$ . This model is described in depth later, in Sec. 20.10.3.

Using approximate inference, it is possible to train and use a wide variety of models. Many of these models are described in the next chapter.

# Chapter 20

# Deep Generative Models

In this chapter, we present several of the specific kinds of generative models that can be built and trained using the techniques presented in Chapters 16, 17, 18 and 19. All of these models represent probability distributions over multiple variables in some way. Some allow the probability distribution function to be evaluated explicitly. Others do not allow the evaluation of the probability distribution function, but support operations that implicitly require knowledge of it, such as drawing samples from the distribution. Some of these models are structured probabilistic models described in terms of graphs and factors, using the language of graphical models presented in Chapter 16. Others can not easily be described in terms of factors, but represent probability distributions nonetheless.

## 20.1 Boltzmann Machines

Boltzmann machines were originally introduced as a general “connectionist” approach to learning arbitrary probability distributions over binary vectors (Fahlman *et al.*, 1983; Ackley *et al.*, 1985; Hinton *et al.*, 1984; Hinton and Sejnowski, 1986). Variants of the Boltzmann machine that include other kinds of variables have long ago surpassed the popularity of the original. In this section we briefly introduce the binary Boltzmann machine and discuss the issues that come up when trying to train and perform inference in the model.

We define the Boltzmann machine over a  $d$ -dimensional binary random vector  $\mathbf{x} \in \{0, 1\}^d$ . The Boltzmann machine is an energy-based model (Sec. 16.2.4),

meaning we define the joint probability distribution using an energy function:

$$P(\mathbf{x}) = \frac{\exp(-E(\mathbf{x}))}{Z}, \quad (20.1)$$

where  $E(\mathbf{x})$  is the energy function and  $Z$  is the partition function that ensures that  $\sum_{\mathbf{x}} P(\mathbf{x}) = 1$ . The energy function of the Boltzmann machine is given by

$$E(\mathbf{x}) = -\mathbf{x}^\top \mathbf{U}\mathbf{x} - \mathbf{b}^\top \mathbf{x}, \quad (20.2)$$

where  $\mathbf{U}$  is the “weight” matrix of model parameters and  $\mathbf{b}$  is the vector of bias parameters.

In the general setting of the Boltzmann machine, we are given a set of training examples, each of which are  $n$ -dimensional. Eq. 20.1 describes the joint probability distribution over the observed variables. While this scenario is certainly viable, it does limit the kinds of interactions between the observed variables to those described by the weight matrix. Specifically, it means that the probability of one unit being on is given by a linear model (logistic regression) from the values of the other units.

The Boltzmann machine becomes more powerful when not all the variables are observed. In this case, the non-observed variables, or *latent* variables, can act similarly to hidden units in a multi-layer perceptron and model higher-order interactions among the visible units. Just as the addition of hidden units to convert logistic regression into an MLP results in the MLP being a universal approximator of functions, a Boltzmann machine with hidden units is no longer limited to modeling linear relationships between variables. Instead, the Boltzmann machine becomes a universal approximator of probability mass functions over discrete variables (Le Roux and Bengio, 2008).

Formally, we decompose the units  $\mathbf{x}$  into two subsets: the visible units  $\mathbf{v}$  and the latent (or hidden) units  $\mathbf{h}$ . The energy function becomes

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{v}^\top \mathbf{R}\mathbf{v} - \mathbf{v}^\top \mathbf{W}\mathbf{h} - \mathbf{h}^\top \mathbf{S}\mathbf{h} - \mathbf{b}^\top \mathbf{v} - \mathbf{c}^\top \mathbf{h}. \quad (20.3)$$

**Boltzmann Machine Learning** Learning algorithms for Boltzmann machines are usually based on maximum likelihood. All Boltzmann machines have an intractable partition function, so the maximum likelihood gradient must be approximated using the techniques described in Chapter 18.

One interesting property of Boltzmann machines when trained with learning rules based on maximum likelihood is that the update for a particular weight connecting two units depends only the statistics of those two units, collected

under different distributions:  $P_{\text{model}}(\mathbf{v})$  and  $\hat{P}_{\text{data}}(\mathbf{v}) P_{\text{model}}(\mathbf{h} \mid \mathbf{v})$ . The rest of the network participates in shaping those statistics, but the weight can be updated without knowing anything about the rest of the network or how those statistics were produced. This means that the learning rule is “local,” which makes Boltzmann machine learning somewhat biologically plausible. It is conceivable that if each neuron were a random variable in a Boltzmann machine, then the axons and dendrites connecting two random variables could learn only by observing the firing pattern of the cells that they actually physically touch. In particular, in the positive phase, two units that frequently activate together have their connection strengthened. This is an example of a Hebbian learning rule (Hebb, 1949) often summarized with the mnemonic “fire together, wire together.” Hebbian learning rules are among the oldest hypothesized explanations for learning in biological systems and remain relevant today (Giudice *et al.*, 2009).

Other learning algorithms that use more information than local statistics seem to require us to hypothesize the existence of more machinery than this. For example, for the brain to implement back-propagation in a multilayer perceptron, it seems necessary for the brain to maintain a secondary communication network for transmitting gradient information backwards through the network. Proposals for biologically plausible implementations (and approximations) of back-propagation have been made (Hinton, 2007a; Bengio, 2015) but remain to be validated, and Bengio (2015) links back-propagation of gradients to inference in energy-based models similar to the Boltzmann machine (but with continuous latent variables).

The negative phase of Boltzmann machine learning is somewhat harder to explain from a biological point of view. As argued in Sec. 18.2, dream sleep may be a form of negative phase sampling. This idea is more speculative though.

## 20.2 Restricted Boltzmann Machines

Invented under the name *harmonium* (Smolensky, 1986), restricted Boltzmann machines are some of the most common building blocks of deep probabilistic models. We have briefly described RBMs previously, in Sec. 16.7.1. Here we review the previous information and go into more detail. RBMs are undirected probabilistic graphical models containing a layer of observable variables and a single layer of latent variables. RBMs may be stacked (one on top of the other) to form deeper models. See Fig. 20.1 for some examples. In particular, Fig. 20.1a shows the graph structure of the RBM itself. It is a bipartite graph, with no connections permitted between any variables in the observed layer or between any units in the latent layer.

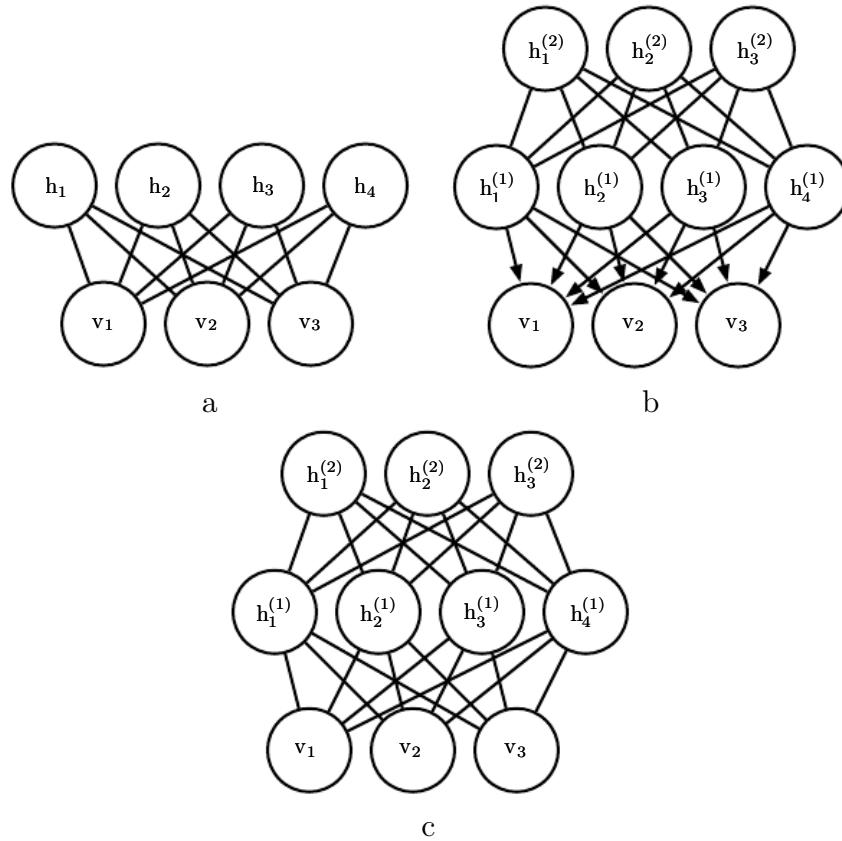


Figure 20.1: Examples of models that may be built with restricted Boltzmann machines. (a) The restricted Boltzmann machine itself is an undirected graphical model based on a bipartite graph, with visible units in one part of the graph and hidden units in the other part. There are no connections among the visible units, nor any connections among the hidden units. Typically every visible unit is connected to every hidden unit but it is possible to construct sparsely connected RBMs such as convolutional RBMs. (b) A deep belief network is a hybrid graphical model involving both directed and undirected connections. Like an RBM, it has no intra-layer connections. However, a DBN has multiple hidden layers, and thus there are connections between hidden units that are in separate layers. All of the local conditional probability distributions needed by the deep belief network are copied directly from the local conditional probability distributions of its constituent RBMs. Alternatively, we could also represent the deep belief network with a completely undirected graph, but it would need intra-layer connections to capture the dependencies between parents. (c) A deep Boltzmann machine is an undirected graphical model with several layers of latent variables. Like RBMs and DBNs, DBMs lack intra-layer connections. DBMs are less closely tied to RBMs than DBNs are. When initializing a DBM from a stack of RBMs, it is necessary to modify the RBM parameters slightly. Some kinds of DBMs may be trained without first training a set of RBMs.

We begin with the binary version of the restricted Boltzmann machine, but as we see later there are extensions to other types of visible and hidden units.

More formally, let the observed layer consist of a set of  $n_v$  binary random variables which we refer to collectively with the vector  $\mathbf{v}$ . We refer to the latent or hidden layer of  $n_h$  binary random variables as  $\mathbf{h}$ .

Like the general Boltzmann machine, the restricted Boltzmann machine is an energy-based model with the joint probability distribution specified by its energy function:

$$P(\mathbf{v} = \mathbf{v}, \mathbf{h} = \mathbf{h}) = \frac{1}{Z} \exp(-E(\mathbf{v}, \mathbf{h})). \quad (20.4)$$

The energy function for an RBM is given by

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{b}^\top \mathbf{v} - \mathbf{c}^\top \mathbf{h} - \mathbf{v}^\top \mathbf{W} \mathbf{h}, \quad (20.5)$$

and  $Z$  is the normalizing constant known as the partition function:

$$Z = \sum_{\mathbf{v}} \sum_{\mathbf{h}} \exp\{-E(\mathbf{v}, \mathbf{h})\}. \quad (20.6)$$

It is apparent from the definition of the partition function  $Z$  that the naive method of computing  $Z$  (exhaustively summing over all states) could be computationally intractable, unless a cleverly designed algorithm could exploit regularities in the probability distribution to compute  $Z$  faster. In the case of restricted Boltzmann machines, [Long and Servedio \(2010\)](#) formally proved that the partition function  $Z$  is intractable. The intractable partition function  $Z$  implies that the normalized joint probability distribution  $P(\mathbf{v})$  is also intractable to evaluate.

### 20.2.1 Conditional Distributions

Though  $P(\mathbf{v})$  is intractable, the bipartite graph structure of the RBM has the very special property that its conditional distributions  $P(\mathbf{h} | \mathbf{v})$  and  $P(\mathbf{v} | \mathbf{h})$  are factorial and relatively simple to compute and to sample from.

Deriving the conditional distributions from the joint distribution is straightforward:

$$P(\mathbf{h} | \mathbf{v}) = \frac{P(\mathbf{h}, \mathbf{v})}{P(\mathbf{v})} \quad (20.7)$$

$$= \frac{1}{P(\mathbf{v})} \frac{1}{Z} \exp \left\{ \mathbf{b}^\top \mathbf{v} + \mathbf{c}^\top \mathbf{h} + \mathbf{v}^\top \mathbf{W} \mathbf{h} \right\} \quad (20.8)$$

$$= \frac{1}{Z'} \exp \left\{ \mathbf{c}^\top \mathbf{h} + \mathbf{v}^\top \mathbf{W} \mathbf{h} \right\} \quad (20.9)$$

$$= \frac{1}{Z'} \exp \left\{ \sum_{j=1}^{n_h} c_j h_j + \sum_{j=1}^{n_h} \mathbf{v}^\top \mathbf{W}_{:,j} \mathbf{h}_j \right\} \quad (20.10)$$

$$= \frac{1}{Z'} \prod_{j=1}^{n_h} \exp \left\{ c_j h_j + \mathbf{v}^\top \mathbf{W}_{:,j} \mathbf{h}_j \right\} \quad (20.11)$$

Since we are conditioning on the visible units  $\mathbf{v}$ , we can treat these as constant with respect to the distribution  $P(\mathbf{h} | \mathbf{v})$ . The factorial nature of the conditional  $P(\mathbf{h} | \mathbf{v})$  follows immediately from our ability to write the joint probability over the vector  $\mathbf{h}$  as the product of (unnormalized) distributions over the individual elements,  $h_j$ . It is now a simple matter of normalizing the distributions over the individual binary  $h_j$ .

$$P(h_j = 1 | \mathbf{v}) = \frac{\tilde{P}(h_j = 1 | \mathbf{v})}{\tilde{P}(h_j = 0 | \mathbf{v}) + \tilde{P}(h_j = 1 | \mathbf{v})} \quad (20.12)$$

$$= \frac{\exp \{ c_j + \mathbf{v}^\top \mathbf{W}_{:,j} \}}{\exp \{ 0 \} + \exp \{ c_j + \mathbf{v}^\top \mathbf{W}_{:,j} \}} \quad (20.13)$$

$$= \sigma(c_j + \mathbf{v}^\top \mathbf{W}_{:,j}). \quad (20.14)$$

We can now express the full conditional over the hidden layer as the factorial distribution:

$$P(\mathbf{h} | \mathbf{v}) = \prod_{j=1}^{n_h} \sigma((2\mathbf{h} - 1) \odot (\mathbf{c} + \mathbf{W}^\top \mathbf{v}))_j. \quad (20.15)$$

A similar derivation will show that the other condition of interest to us,  $P(\mathbf{v} | \mathbf{h})$ , is also a factorial distribution:

$$P(\mathbf{v} | \mathbf{h}) = \prod_{i=1}^{n_v} \sigma((2\mathbf{v} - 1) \odot (\mathbf{b} + \mathbf{W}\mathbf{h}))_i. \quad (20.16)$$

### 20.2.2 Training Restricted Boltzmann Machines

Because the RBM admits efficient evaluation and differentiation of  $\tilde{P}(\mathbf{v})$  and efficient MCMC sampling in the form of block Gibbs sampling, it can readily be trained with any of the techniques described in Chapter 18 for training models that have intractable partition functions. This includes CD, SML (PCD), ratio matching and so on. Compared to other undirected models used in deep learning, the RBM is relatively straightforward to train because we can compute  $P(\mathbf{h} | \mathbf{v})$

exactly in closed form. Some other deep models, such as the deep Boltzmann machine, combine both the difficulty of an intractable partition function and the difficulty of intractable inference.

## 20.3 Deep Belief Networks

*Deep belief networks* (DBNs) were one of the first non-convolutional models to successfully admit training of deep architectures (Hinton *et al.*, 2006; Hinton, 2007b). The introduction of deep belief networks in 2006 began the current deep learning renaissance. Prior to the introduction of deep belief networks, deep models were considered too difficult to optimize. Kernel machines with convex objective functions dominated the research landscape. Deep belief networks demonstrated that deep architectures can be successful, by outperforming kernelized support vector machines on the MNIST dataset (Hinton *et al.*, 2006). Today, deep belief networks have mostly fallen out of favor and are rarely used, even compared to other unsupervised or generative learning algorithms, but they are still deservedly recognized for their important role in deep learning history.

Deep belief networks are generative models with several layers of latent variables. The latent variables are typically binary, while the visible units may be binary or real. There are no intra-layer connections. Usually, every unit in each layer is connected to every unit in each neighboring layer, though it is possible to construct more sparsely connected DBNs. The connections between the top two layers are undirected. The connections between all other layers are directed, with the arrows pointed toward the layer that is closest to the data. See Fig. 20.1b for an example.

A DBN with  $l$  hidden layers contains  $l$  weight matrices:  $\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(l)}$ . It also contains  $l+1$  bias vectors:  $\mathbf{b}^{(0)}, \dots, \mathbf{b}^{(l)}$ , with  $\mathbf{b}^{(0)}$  providing the biases for the visible layer. The probability distribution represented by the DBN is given by

$$P(\mathbf{h}^{(l)}, \mathbf{h}^{(l-1)}) \propto \exp\left(\mathbf{b}^{(l)\top} \mathbf{h}^{(l)} + \mathbf{b}^{(l-1)\top} \mathbf{h}^{(l-1)} + \mathbf{h}^{(l-1)\top} \mathbf{W}^{(l)} \mathbf{h}^{(l)}\right), \quad (20.17)$$

$$P(h_i^{(k)} = 1 | \mathbf{h}^{(k+1)}) = \sigma\left(b_i^{(k)} + \mathbf{W}_{:,i}^{(k+1)\top} \mathbf{h}^{(k+1)}\right) \forall i, \forall k \in 1, \dots, l-2, \quad (20.18)$$

$$P(v_i = 1 | \mathbf{h}^{(1)}) = \sigma\left(b_i^{(0)} + \mathbf{W}_{:,i}^{(1)\top} \mathbf{h}^{(1)}\right) \forall i. \quad (20.19)$$

In the case of real-valued visible units, substitute

$$\mathbf{v} \sim \mathcal{N}\left(\mathbf{v}; \mathbf{b}^{(0)} + \mathbf{W}^{(1)\top} \mathbf{h}^{(1)}, \boldsymbol{\beta}^{-1}\right) \quad (20.20)$$

with  $\beta$  diagonal for tractability. Generalizations to other exponential family visible units are straightforward, at least in theory. A DBN with only one hidden layer is just an RBM.

To generate a sample from a DBN, we first run several steps of Gibbs sampling on the top two hidden layers. This stage is essentially drawing a sample from the RBM defined by the top two hidden layers. We can then use a single pass of ancestral sampling through the rest of the model to draw a sample from the visible units.

Deep belief networks incur many of the problems associated with both directed models and undirected models.

Inference in a deep belief network is intractable due to the explaining away effect within each directed layer, and due to the interaction between the two hidden layers that have undirected connections. Evaluating or maximizing the standard evidence lower bound on the log-likelihood is also intractable, because the evidence lower bound takes the expectation of cliques whose size is equal to the network width.

Evaluating or maximizing the log-likelihood requires not just confronting the problem of intractable inference to marginalize out the latent variables, but also the problem of an intractable partition function within the undirected model of the top two layers.

To train a deep belief network, one begins by training an RBM to maximize  $\mathbb{E}_{\mathbf{v} \sim p_{\text{data}}} \log p(\mathbf{v})$  using contrastive divergence or stochastic maximum likelihood. The parameters of the RBM then define the parameters of the first layer of the DBN. Next, a second RBM is trained to approximately maximize

$$\mathbb{E}_{\mathbf{v} \sim p_{\text{data}}} \mathbb{E}_{\mathbf{h}^{(1)} \sim p^{(1)}(\mathbf{h}^{(1)} | \mathbf{v})} \log p^{(2)}(\mathbf{h}^{(1)}) \quad (20.21)$$

where  $p^{(1)}$  is the probability distribution represented by the first RBM and  $p^{(2)}$  is the probability distribution represented by the second RBM. In other words, the second RBM is trained to model the distribution defined by sampling the hidden units of the first RBM, when the first RBM is driven by the data. This procedure can be repeated indefinitely, to add as many layers to the DBN as desired, with each new RBM modeling the samples of the previous one. Each RBM defines another layer of the DBN. This procedure can be justified as increasing a variational lower bound on the log-likelihood of the data under the DBN ([Hinton et al., 2006](#)).

In most applications, no effort is made to jointly train the DBN after the greedy layer-wise procedure is complete. However, it is possible to perform generative fine-tuning using the wake-sleep algorithm.

The trained DBN may be used directly as a generative model, but most of the interest in DBNs arose from their ability to improve classification models. We can take the weights from the DBN and use them to define an MLP:

$$\mathbf{h}^{(1)} = \sigma \left( b^{(1)} + \mathbf{v}^\top \mathbf{W}^{(1)} \right). \quad (20.22)$$

$$\mathbf{h}^{(l)} = \sigma \left( b_i^{(l)} + \mathbf{h}^{(l-1)\top} \mathbf{W}^{(l)} \right) \forall l \in 2, \dots, m, \quad (20.23)$$

After initializing this MLP with the weights and biases learned via generative training of the DBN, we may train the MLP to perform a classification task. This additional training of the MLP is an example of discriminative fine-tuning.

This specific choice of MLP is somewhat arbitrary, compared to many of the inference equations in Chapter 19 that are derived from first principles. This MLP is a heuristic choice that seems to work well in practice and is used consistently in the literature. Many approximate inference techniques are motivated by their ability to find a maximally *tight* variational lower bound on the log-likelihood under some set of constraints. One can construct a variational lower bound on the log-likelihood using the hidden unit expectations defined by the DBN’s MLP, but this is true of *any* probability distribution over the hidden units, and there is no reason to believe that this MLP provides a particularly tight bound. In particular, the MLP ignores many important interactions in the DBN graphical model. The MLP propagates information upward from the visible units to the deepest hidden units, but does not propagate any information downward or sideways. The DBN graphical model has explaining away interactions between all of the hidden units within the same layer as well as top-down interactions between layers.

While the log-likelihood of a DBN is intractable, it may be approximated with AIS (Salakhutdinov and Murray, 2008). This permits evaluating its quality as a generative model.

The term “deep belief network” is commonly used incorrectly to refer to any kind of deep neural network, even networks without latent variable semantics. The term “deep belief network” should refer specifically to models with undirected connections in the deepest layer and directed connections pointing downward between all other pairs of consecutive layers.

The term “deep belief network” may also cause some confusion because the term “belief network” is sometimes used to refer to purely directed models, while deep belief networks contain an undirected layer. Deep belief networks also share the acronym DBN with dynamic Bayesian networks (Dean and Kanazawa, 1989), which are Bayesian networks for representing Markov chains.

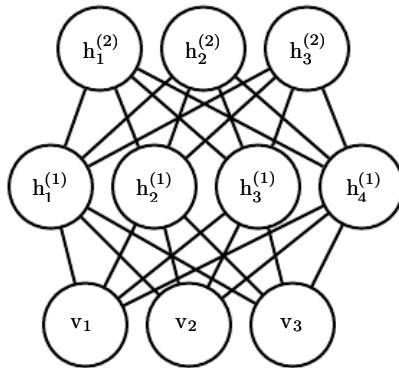


Figure 20.2: The graphical model for a deep Boltzmann machine with one visible layer (bottom) and two hidden layers. Connections are only between units in neighboring layers. There are no intra-layer layer connections.

## 20.4 Deep Boltzmann Machines

A *deep Boltzmann machine* or *DBM* (Salakhutdinov and Hinton, 2009a) is another kind of deep, generative model. Unlike the deep belief network (DBN), it is an entirely undirected model. Unlike the RBM, the DBM has several layers of latent variables (RBMs have just one). But like the RBM, within each layer, each of the variables are mutually independent, conditioned on the variables in the neighboring layers. See Fig. 20.2 for the graph structure. Deep Boltzmann machines have been applied to a variety of tasks including document modeling (Srivastava *et al.*, 2013).

Like RBMs and DBNs, DBMs typically contain only binary units—as we assume for simplicity of our presentation of the model—but it is straightforward to include real-valued visible units.

A DBM is an energy-based model, meaning that the joint probability distribution over the model variables is parametrized by an energy function  $E$ . In the case of a deep Boltzmann machine with one visible layer,  $\mathbf{v}$ , and three hidden layers,  $\mathbf{h}^{(1)}$ ,  $\mathbf{h}^{(2)}$  and  $\mathbf{h}^{(3)}$ , the joint probability is given by:

$$P(\mathbf{v}, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}, \mathbf{h}^{(3)}) = \frac{1}{Z(\boldsymbol{\theta})} \exp(-E(\mathbf{v}, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}, \mathbf{h}^{(3)}; \boldsymbol{\theta})). \quad (20.24)$$

To simplify our presentation, we omit the bias parameters below. The DBM energy function is then defined as follows:

$$E(\mathbf{v}, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}, \mathbf{h}^{(3)}; \boldsymbol{\theta}) = -\mathbf{v}^\top \mathbf{W}^{(1)} \mathbf{h}^{(1)} - \mathbf{h}^{(1)\top} \mathbf{W}^{(2)} \mathbf{h}^{(2)} - \mathbf{h}^{(2)\top} \mathbf{W}^{(3)} \mathbf{h}^{(3)}. \quad (20.25)$$

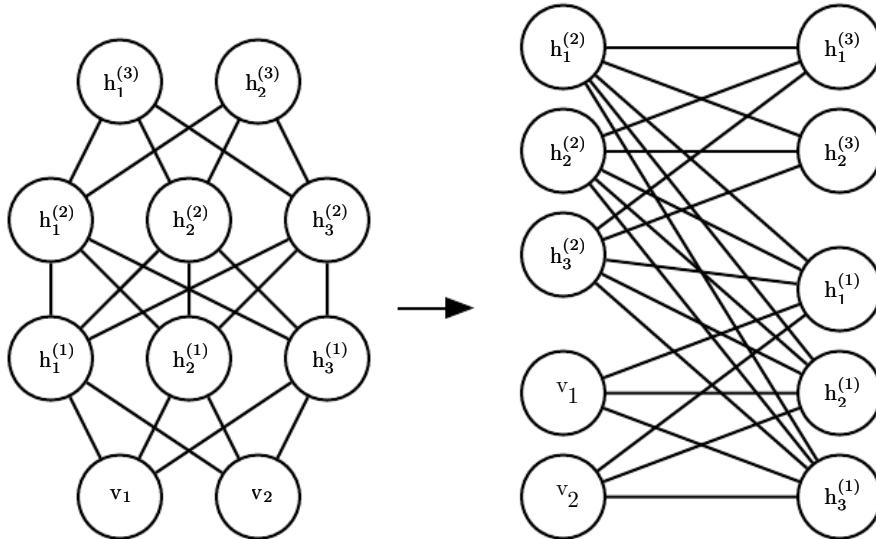


Figure 20.3: A deep Boltzmann machine, re-arranged to reveal its bipartite graph structure.

In comparison to the RBM energy function (Eq. 20.5), the DBM energy function includes connections between the hidden units (latent variables) in the form of the weight matrices ( $\mathbf{W}^{(2)}$  and  $\mathbf{W}^{(3)}$ ). As we will see, these connections have significant consequences for both the model behavior as well as how we go about performing inference in the model.

In comparison to fully connected Boltzmann machines (with every unit connected to every other unit), the DBM offers some advantages that are similar to those offered by the RBM. Specifically, as illustrated in Fig. 20.3, the DBM layers can be organized into a bipartite graph, with odd layers on one side and even layers on the other. This immediately implies that when we condition on the variables in the even layer, the variables in the odd layers become conditionally independent. Of course, when we condition on the variables in the odd layers, the variables in the even layers also become conditionally independent.

The bipartite structure of the DBM means that we can apply the same equations we have previously used for the conditional distributions of an RBM to determine the conditional distributions in a DBM. The units within a layer are conditionally independent from each other given the values of the neighboring layers, so the distributions over binary variables can be fully described by the Bernoulli parameters giving the probability of each unit being active. In our example with two hidden layers, the activation probabilities are given by:

$$P(v_i = 1 \mid \mathbf{h}^{(1)}) = \sigma(\mathbf{W}_{i,:}^{(1)} \mathbf{h}^{(1)}), \quad (20.26)$$

$$P(h_i^{(1)} = 1 \mid \mathbf{v}, \mathbf{h}^{(2)}) = \sigma(\mathbf{v}^\top \mathbf{W}_{:,i}^{(1)} + \mathbf{W}_{i,:}^{(2)} \mathbf{h}^{(2)}) \quad (20.27)$$

and

$$P(h_k^{(2)} = 1 \mid \mathbf{h}^{(1)}) = \sigma(\mathbf{h}^{(1)\top} \mathbf{W}_{:,k}^{(2)}). \quad (20.28)$$

The bipartite structure makes Gibbs sampling in a deep Boltzmann machine efficient. The naive approach to Gibbs sampling is to update only one variable at a time. RBMs allow all of the visible units to be updated in one block and all of the hidden units to be updated in a second block. One might naively assume that a DBM with  $l$  layers requires  $l+1$  updates, with each iteration updating a block consisting of one layer of units. Instead, it is possible to update all of the units in only two iterations. Gibbs sampling can be divided into two blocks of updates, one including all even layers (including the visible layer) and the other including all odd layers. Due to the bipartite DBM connection pattern, given the even layers, the distribution over the odd layers is factorial and thus can be sampled simultaneously and independently as a block. Likewise, given the odd layers, the even layers can be sampled simultaneously and independently as a block. Efficient sampling is especially important for training with the stochastic maximum likelihood algorithm.

### 20.4.1 Interesting Properties

Deep Boltzmann machines have many interesting properties.

DBMs were developed after DBNs. Compared to DBNs, the posterior distribution  $P(\mathbf{h} \mid \mathbf{v})$  is simpler for DBMs. Somewhat counterintuitively, the simplicity of this posterior distribution allows richer approximations of the posterior. In the case of the DBN, we perform classification using a heuristically motivated approximate inference procedure, in which we guess that a reasonable value for the mean field expectation of the hidden units can be provided by an upward pass through the network in an MLP that uses sigmoid activation functions and the same weights as the original DBN. **Any** distribution  $Q(\mathbf{h})$  may be used to obtain a variational lower bound on the log-likelihood. This heuristic procedure therefore allows us to obtain such a bound. However, the bound is not explicitly optimized in any way, so the bound may be far from tight. In particular, the heuristic estimate of  $Q$  ignores interactions between hidden units within the same layer as well as the top-down feedback influence of hidden units in deeper layers on hidden units that are closer to the input. Because the heuristic MLP-based inference procedure in the DBN is not able to account for these interactions, the resulting  $Q$  is presumably far

from optimal. In DBMs, all of the hidden units within a layer are conditionally independent given the other layers. This lack of intra-layer interaction makes it possible to use fixed point equations to actually optimize the variational lower bound and find the true optimal mean field expectations (to within some numerical tolerance).

The use of proper mean field allows the approximate inference procedure for DBMs to capture the influence of top-down feedback interactions. This makes DBMs interesting from the point of view of neuroscience, because the human brain is known to use many top-down feedback connections. Because of this property, DBMs have been used as computational models of real neuroscientific phenomena (Series *et al.*, 2010; Reichert *et al.*, 2011).

One unfortunate property of DBMs is that sampling from them is relatively difficult. DBNs only need to use MCMC sampling in their top pair of layers. The other layers are used only at the end of the sampling process, in one efficient ancestral sampling pass. To generate a sample from a DBM, it is necessary to use MCMC across all layers, with every layer of the model participating in every Markov chain transition.

#### 20.4.2 DBM Mean Field Inference

The conditional distribution over one DBM layer given the neighboring layers is factorial. In the example of the DBM with two hidden layers, these distributions are  $P(\mathbf{v} \mid \mathbf{h}^{(1)})$ ,  $P(\mathbf{h}^{(1)} \mid \mathbf{v}, \mathbf{h}^{(2)})$  and  $P(\mathbf{h}^{(2)} \mid \mathbf{h}^{(1)})$ . The distribution over **all** hidden layers generally does not factorize because of interactions between layers. In the example with two hidden layers,  $P(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} \mid \mathbf{v})$  does not factorize due to the interaction weights  $\mathbf{W}^{(2)}$  between  $\mathbf{h}^{(1)}$  and  $\mathbf{h}^{(2)}$  which render these variables mutually dependent.

As was the case with the DBN, we are left to seek out methods to approximate the DBM posterior distribution. However, unlike the DBN, the DBM posterior distribution over their hidden units—while complicated—is easy to approximate with a *variational* approximation (as discussed in Sec. 19.4), specifically a mean field approximation. The mean field approximation is a simple form of variational inference, where we restrict the approximating distribution to fully factorial distributions. In the context of DBMs, the mean field equations capture the bidirectional interactions between layers. In this section we derive the iterative approximate inference procedure originally introduced in Salakhutdinov and Hinton (2009a).

In variational approximations to inference, we approach the task of approximating a particular target distribution—in our case, the posterior distribution over

the hidden units given the visible units—by some reasonably simple family of distributions. In the case of the mean field approximation, the approximating family is the set of distributions where the hidden units are conditionally independent.

We now develop the mean field approach for the example with two hidden layers. Let  $Q(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v})$  be the approximation of  $P(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v})$ . The mean field assumption implies that

$$Q(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v}) = \prod_j Q(h_j^{(1)} | \mathbf{v}) \prod_k Q(h_k^{(2)} | \mathbf{v}). \quad (20.29)$$

The mean field approximation attempts to find a member of this family of distributions that best fits the true posterior  $P(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v})$ . Importantly, the inference process must be run again to find a different distribution  $Q$  every time we use a new value of  $\mathbf{v}$ .

One can conceive of many ways of measuring how well  $Q(\mathbf{h} | \mathbf{v})$  fits  $P(\mathbf{h} | \mathbf{v})$ . The mean field approach is to minimize

$$\text{KL}(Q \| P) = \sum_{\mathbf{h}} Q(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v}) \log \frac{Q(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v})}{P(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v})}. \quad (20.30)$$

In general, we do not have to provide a parametric form of the approximating distribution beyond enforcing the independence assumptions. The variational approximation procedure is generally able to recover a functional form of the approximate distribution. However, in the case of a mean field assumption on binary hidden units (the case we are developing here) there is no loss of generality resulting from fixing a parametrization of the model in advance.

We parametrize  $Q$  as a product of Bernoulli distributions, that is we associate the probability of each element of  $\mathbf{h}^{(1)}$  with a parameter. Specifically, for each  $j$ ,  $\hat{h}_j^{(1)} = Q(h_j^{(1)} = 1 | \mathbf{v})$ , where  $\hat{h}_j^{(1)} \in [0, 1]$  and for each  $k$ ,  $\hat{h}_k^{(2)} = Q(h_k^{(2)} = 1 | \mathbf{v})$ , where  $\hat{h}_k^{(2)} \in [0, 1]$ . Thus we have the following approximation to the posterior:

$$Q(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v}) = \prod_j Q(h_j^{(1)} | \mathbf{v}) \prod_k Q(h_k^{(2)} | \mathbf{v}) \quad (20.31)$$

$$= \prod_j (\hat{h}_j^{(1)})^{h_j^{(1)}} (1 - \hat{h}_j^{(1)})^{(1-h_j^{(1)})} \times \prod_k (\hat{h}_k^{(2)})^{h_k^{(2)}} (1 - \hat{h}_k^{(2)})^{(1-h_k^{(2)})}. \quad (20.32)$$

Of course, for DBMs with more layers the approximate posterior parametrization can be extended in the obvious way, exploiting the bipartite structure of the graph

to update all of the even layers simultaneously and then to update all of the odd layers simultaneously, following the same schedule as Gibbs sampling.

Now that we have specified our family of approximating distributions  $Q$ , it remains to specify a procedure for choosing the member of this family that best fits  $P$ . The most straightforward way to do this is to use the mean field equations specified by Eq. 19.56. These equations were derived by solving for where the derivatives of the variational lower bound are zero. They describe in an abstract manner how to optimize the variational lower bound for any model, simply by taking expectations with respect to  $Q$ .

Applying these general equations, we obtain the update rules (again, ignoring bias terms):

$$\hat{h}_j^{(1)} = \sigma \left( \sum_i v_i W_{i,j}^{(1)} + \sum_{k'} W_{j,k'}^{(2)} \hat{h}_{k'}^{(2)} \right), \quad \forall j \quad (20.33)$$

$$\hat{h}_k^{(2)} = \sigma \left( \sum_{j'} W_{j',k}^{(2)} \hat{h}_{j'}^{(1)} \right), \quad \forall k. \quad (20.34)$$

At a fixed point of this system of equations, we have a local maximum of the variational lower bound  $\mathcal{L}(Q)$ . Thus these fixed point update equations define an iterative algorithm where we alternate updates of  $\hat{h}_j^{(1)}$  (using Eq. 20.33) and updates of  $\hat{h}_k^{(2)}$  (using Eq. 20.34). On small problems such as MNIST, as few as ten iterations can be sufficient to find an approximate positive phase gradient for learning, and fifty usually suffice to obtain a high quality representation of a single specific example to be used for high-accuracy classification. Extending approximate variational inference to deeper DBMs is straightforward.

### 20.4.3 DBM Parameter Learning

Learning in the DBM must confront both the challenge of an intractable partition function, using the techniques from Chapter 18, and the challenge of an intractable posterior distribution, using the techniques from Chapter 19.

As described in Sec. 20.4.2, variational inference allows the construction of a distribution  $Q(\mathbf{h} | \mathbf{v})$  that approximates the intractable  $P(\mathbf{h} | \mathbf{v})$ . Learning then proceeds by maximizing  $\mathcal{L}(\mathbf{v}, Q, \boldsymbol{\theta})$ , the variational lower bound on the intractable log-likelihood,  $\log P(\mathbf{v}; \boldsymbol{\theta})$ .

For a deep Boltzmann machine with two hidden layers,  $\mathcal{L}$  is given by

$$\mathcal{L}(Q, \boldsymbol{\theta}) = \sum_i \sum_{j'} v_i W_{i,j'}^{(1)} \hat{h}_{j'}^{(1)} + \sum_{j'} \sum_{k'} \hat{h}_{j'}^{(1)} W_{j',k'}^{(2)} \hat{h}_{k'}^{(2)} - \log Z(\boldsymbol{\theta}) + \mathcal{H}(Q). \quad (20.35)$$

This expression still contains the log partition function,  $\log Z(\boldsymbol{\theta})$ . Because a deep Boltzmann machine contains restricted Boltzmann machines as components, the hardness results for computing the partition function and sampling that apply to restricted Boltzmann machines also apply to deep Boltzmann machines. This means that evaluating the probability mass function of a Boltzmann machine requires approximate methods such as annealed importance sampling. Likewise, training the model requires approximations to the gradient of the log partition function. See Chapter 18 for a general description of these methods. DBMs are typically trained using stochastic maximum likelihood. Many of the other techniques described in Chapter 18 are not applicable. Techniques such as pseudolikelihood require the ability to evaluate the unnormalized probabilities, rather than merely obtain a variational lower bound on them. Contrastive divergence is slow for deep Boltzmann machines because they do not allow efficient sampling of the hidden units given the visible units—instead, contrastive divergence would require burning in a Markov chain every time a new negative phase sample is needed.

The non-variational version of stochastic maximum likelihood algorithm was discussed earlier, in Sec. 18.2. Variational stochastic maximum likelihood as applied to the DBM is given in Algorithm 20.1. Recall that we describe a simplified variant of the DBM that lacks bias parameters; including them is trivial.

#### 20.4.4 Layer-Wise Pretraining

Unfortunately, training a DBM using stochastic maximum likelihood (as described above) from a random initialization usually results in failure. In some cases, the model fails to learn to represent the distribution adequately. In other cases, the DBM may represent the distribution well, but with no higher likelihood than could be obtained with just an RBM. A DBM with very small weights in all but the first layer represents approximately the same distribution as an RBM.

Various techniques that permit joint training have been developed and are described in Sec. 20.4.5. However, the original and most popular method for overcoming the joint training problem of DBMs is greedy layer-wise pretraining. In this method, each layer of the DBM is trained in isolation as an RBM. The first layer is trained to model the input data. Each subsequent RBM is trained to model samples from the previous RBM’s posterior distribution. After all of the

---

**Algorithm 20.1** The variational stochastic maximum likelihood algorithm for training a DBM with two hidden layers.

---

Set  $\epsilon$ , the step size, to a small positive number  
Set  $k$ , the number of Gibbs steps, high enough to allow a Markov chain of  $p(\mathbf{v}, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}; \boldsymbol{\theta} + \epsilon \Delta_{\boldsymbol{\theta}})$  to burn in, starting from samples from  $p(\mathbf{v}, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}; \boldsymbol{\theta})$ . Initialize three matrices,  $\tilde{\mathbf{V}}$ ,  $\tilde{\mathbf{H}}^{(1)}$  and  $\tilde{\mathbf{H}}^{(2)}$  each with  $m$  columns set to random values (e.g., from Bernoulli distributions, possibly with marginals matched to the model's marginals).  
**while** not converged (learning loop) **do**  
    Sample a minibatch of  $m$  examples from the training data and arrange them as the rows of a design matrix  $\mathbf{V}$ .  
    Initialize matrices  $\hat{\mathbf{H}}^{(1)}$  and  $\hat{\mathbf{H}}^{(2)}$ , possibly to the model's marginals.  
    **while** not converged (mean field inference loop) **do**  
         $\hat{\mathbf{H}}^{(1)} \leftarrow \sigma(\mathbf{V}\mathbf{W}^{(1)} + \hat{\mathbf{H}}^{(2)}\mathbf{W}^{(2)\top})$ .  
         $\hat{\mathbf{H}}^{(2)} \leftarrow \sigma(\hat{\mathbf{H}}^{(1)}\mathbf{W}^{(2)})$ .  
    **end while**  
     $\Delta_{\mathbf{W}^{(1)}} \leftarrow \frac{1}{m} \mathbf{V}^\top \hat{\mathbf{H}}^{(1)}$   
     $\Delta_{\mathbf{W}^{(2)}} \leftarrow \frac{1}{m} \hat{\mathbf{H}}^{(1)\top} \hat{\mathbf{H}}^{(2)}$   
    **for**  $l = 1$  to  $k$  (Gibbs sampling) **do**  
        Gibbs block 1:  
         $\forall i, j, \tilde{V}_{i,j}$  sampled from  $P(\tilde{V}_{i,j} = 1) = \sigma(\mathbf{W}_{j,:}^{(1)} \tilde{\mathbf{H}}_{i,:}^{(1)})$ .  
         $\forall i, j, \tilde{H}_{i,j}^{(2)}$  sampled from  $P(\tilde{H}_{i,j}^{(2)} = 1) = \sigma(\tilde{\mathbf{H}}_{i,:}^{(1)} \mathbf{W}_{:,j}^{(2)})$ .  
        Gibbs block 2:  
         $\forall i, j, \tilde{H}_{i,j}^{(1)}$  sampled from  $P(\tilde{H}_{i,j}^{(1)} = 1) = \sigma(\tilde{\mathbf{V}}_{i,:} \mathbf{W}_{:,j}^{(1)} + \tilde{\mathbf{H}}_{i,:}^{(2)} \mathbf{W}_{j,:}^{(2)\top})$ .  
    **end for**  
     $\Delta_{\mathbf{W}^{(1)}} \leftarrow \Delta_{\mathbf{W}^{(1)}} - \frac{1}{m} \mathbf{V}^\top \tilde{\mathbf{H}}^{(1)}$   
     $\Delta_{\mathbf{W}^{(2)}} \leftarrow \Delta_{\mathbf{W}^{(2)}} - \frac{1}{m} \tilde{\mathbf{H}}^{(1)\top} \tilde{\mathbf{H}}^{(2)}$   
     $\mathbf{W}^{(1)} \leftarrow \mathbf{W}^{(1)} + \epsilon \Delta_{\mathbf{W}^{(1)}}$  (this is a cartoon illustration, in practice use a more effective algorithm, such as momentum with a decaying learning rate)  
     $\mathbf{W}^{(2)} \leftarrow \mathbf{W}^{(2)} + \epsilon \Delta_{\mathbf{W}^{(2)}}$   
**end while**

---

RBM<sup>s</sup> have been trained in this way, they can be combined to form a DBM. The DBM may then be trained with PCD. Typically PCD training will make only a small change in the model’s parameters and its performance as measured by the log-likelihood it assigns to the data, or its ability to classify inputs. See Fig. 20.4 for an illustration of the training procedure.

This greedy layer-wise training procedure is not just coordinate ascent. It bears some passing resemblance to coordinate ascent because we optimize one subset of the parameters at each step. However, in the case of the greedy layer-wise training procedure, we actually use a different objective function at each step.

Greedy layer-wise pretraining of a DBM differs from greedy layer-wise pre-training of a DBN. The parameters of each individual RBM may be copied to the corresponding DBN directly. In the case of the DBM, the RBM parameters must be modified before inclusion in the DBM. A layer in the middle of the stack of RBMs is trained with only bottom-up input, but after the stack is combined to form the DBM, the layer will have both bottom-up and top-down input. To account for this effect, Salakhutdinov and Hinton (2009a) advocate dividing the weights of all but the top and bottom RBM in half before inserting them into the DBM. Additionally, the bottom RBM must be trained using two “copies” of each visible unit and the weights tied to be equal between the two copies. This means that the weights are effectively doubled during the upward pass. Similarly, the top RBM should be trained with two copies of the topmost layer.

Obtaining the state of the art results with the deep Boltzmann machine requires a modification of the standard SML algorithm, which is to use a small amount of mean field during the negative phase of the joint PCD training step (Salakhutdinov and Hinton, 2009a). Specifically, the expectation of the energy gradient should be computed with respect to the mean field distribution in which all of the units are independent from each other. The parameters of this mean field distribution should be obtained by running the mean field fixed point equations for just one step. See Goodfellow *et al.* (2013b) for a comparison of the performance of centered DBMs with and without the use of partial mean field in the negative phase.

#### 20.4.5 Jointly Training Deep Boltzmann Machines

Classic DBMs require greedy unsupervised pretraining, and to perform classification well, require a separate MLP-based classifier on top of the hidden features they extract. This has some undesirable properties. It is hard to track performance during training because we cannot evaluate properties of the full DBM while training the first RBM. Thus, it is hard to tell how well our hyperparameters

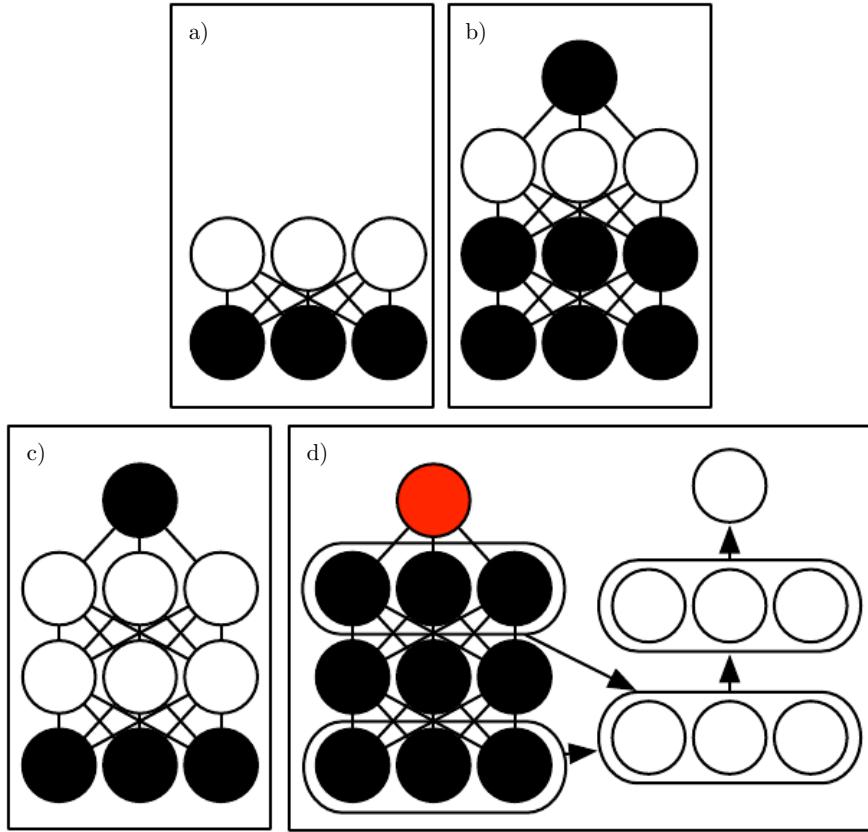


Figure 20.4: The deep Boltzmann machine training procedure used to classify the MNIST dataset (Salakhutdinov and Hinton, 2009a; Srivastava *et al.*, 2014). (a) Train an RBM by using CD to approximately maximize  $\log P(\mathbf{v})$ . (b) Train a second RBM that models  $\mathbf{h}^{(1)}$  and target class  $y$  by using CD- $k$  to approximately maximize  $\log P(\mathbf{h}^{(1)}, y)$  where  $\mathbf{h}^{(1)}$  is drawn from the first RBM's posterior conditioned on the data. Increase  $k$  from 1 to 20 during learning. (c) Combine the two RBMs into a DBM. Train it to approximately maximize  $\log P(\mathbf{v}, y)$  using stochastic maximum likelihood with  $k = 5$ . (d) Delete  $y$  from the model. Define a new set of features  $\mathbf{h}^{(1)}$  and  $\mathbf{h}^{(2)}$  that are obtained by running mean field inference in the model lacking  $y$ . Use these features as input to an MLP whose structure is the same as an additional pass of mean field, with an additional output layer for the estimate of  $y$ . Initialize the MLP's weights to be the same as the DBM's weights. Train the MLP to approximately maximize  $\log P(y | \mathbf{v})$  using stochastic gradient descent and dropout. Figure reprinted from (Goodfellow *et al.*, 2013b).

are working until quite late in the training process. Software implementations of DBMs need to have many different components for CD training of individual RBMs, PCD training of the full DBM, and training based on back-propagation through the MLP. Finally, the MLP on top of the Boltzmann machine loses many of the advantages of the Boltzmann machine probabilistic model, such as being able to perform inference when some input values are missing.

There are two main ways to resolve the joint training problem of the deep Boltzmann machine. The first is the *centered deep Boltzmann machine* (Montavon and Muller, 2012), which reparametrizes the model in order to make the Hessian of the cost function better-conditioned at the beginning of the learning process. This yields a model that can be trained without a greedy layer-wise pretraining stage. The resulting model obtains excellent test set log-likelihood and produces high quality samples. Unfortunately, it remains unable to compete with appropriately regularized MLPs as a classifier. The second way to jointly train a deep Boltzmann machine is to use a *multi-prediction deep Boltzmann machine* (Goodfellow *et al.*, 2013b). This model uses an alternative training criterion that allows the use of the back-propagation algorithm in order to avoid the problems with MCMC estimates of the gradient. Unfortunately, the new criterion does not lead to good likelihood or samples, but, compared to the MCMC approach, it does lead to superior classification performance and ability to reason well about missing inputs.

The centering trick for the Boltzmann machine is easiest to describe if we return to the general view of a Boltzmann machine as consisting of a set of units  $\mathbf{x}$  with a weight matrix  $\mathbf{U}$  and biases  $\mathbf{b}$ . Recall from Eq. 20.2 that the energy function is given by

$$E(\mathbf{x}) = -\mathbf{x}^\top \mathbf{U} \mathbf{x} - \mathbf{b}^\top \mathbf{x}. \quad (20.36)$$

Using different sparsity patterns in the weight matrix  $\mathbf{U}$ , we can implement structures of Boltzmann machines, such as RBMs, or DBMs with different numbers of layers. This is accomplished by partitioning  $\mathbf{x}$  into visible and hidden units and zeroing out elements of  $\mathbf{U}$  for units that do not interact. The centered Boltzmann machine introduces a vector  $\boldsymbol{\mu}$  that is subtracted from all of the states:

$$E'(\mathbf{x}; \mathbf{U}, \mathbf{b}) = -(\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{U} (\mathbf{x} - \boldsymbol{\mu}) - (\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{b}. \quad (20.37)$$

Typically  $\boldsymbol{\mu}$  is a hyperparameter fixed at the beginning of training. It is usually chosen to make sure that  $\mathbf{x} - \boldsymbol{\mu} \approx \mathbf{0}$  when the model is initialized. This reparametrization does not change the set of probability distributions that the model can represent, but it does change the dynamics of stochastic gradient descent applied to the likelihood. Specifically, in many cases, this reparametrization results in a Hessian matrix that is better conditioned. Melchior *et al.* (2013) experimentally

confirmed that the conditioning of the Hessian matrix improves, and observed that the centering trick is equivalent to another Boltzmann machine learning technique, the *enhanced gradient* (Cho *et al.*, 2011). The improved conditioning of the Hessian matrix allows learning to succeed, even in difficult cases like training a deep Boltzmann machine with multiple layers.

The other approach to jointly training deep Boltzmann machines is the multi-prediction deep Boltzmann machine (MP-DBM) which works by viewing the mean field equations as defining a family of recurrent networks for approximately solving every possible inference problem (Goodfellow *et al.*, 2013b). Rather than training the model to maximize the likelihood, the model is trained to make each recurrent network obtain an accurate answer to the corresponding inference problem. The training process is illustrated in Fig. 20.5. It consists of randomly sampling a training example, randomly sampling a subset of inputs to the inference network, and then training the inference network to predict the values of the remaining units.

This general principle of back-propagating through the computational graph for approximate inference has been applied to other models (Stoyanov *et al.*, 2011; Brakel *et al.*, 2013). In these models and in the MP-DBM, the final loss is not the lower bound on the likelihood. Instead, the final loss is typically based on the approximate conditional distribution that the approximate inference network imposes over the missing values. This means that the training of these models is somewhat heuristically motivated. If we inspect the  $p(\mathbf{v})$  represented by the Boltzmann machine learned by the MP-DBM, it tends to be somewhat defective, in the sense that Gibbs sampling yields poor samples.

Back-propagation through the inference graph has two main advantages. First, it trains the model as it is really used—with approximate inference. This means that approximate inference, for example, to fill in missing inputs, or to perform classification despite the presence of missing inputs, is more accurate in the MP-DBM than in the original DBM. The original DBM does not make an accurate classifier on its own; the best classification results with the original DBM were based on training a separate classifier to use features extracted by the DBM, rather than by using inference in the DBM to compute the distribution over the class labels. Mean field inference in the MP-DBM performs well as a classifier without special modifications. The other advantage of back-propagating through approximate inference is that back-propagation computes the exact gradient of the loss. This is better for optimization than the approximate gradients of SML training, which suffer from both bias and variance. This probably explains why MP-DBMs may be trained jointly while DBMs require a greedy layer-wise pretraining.

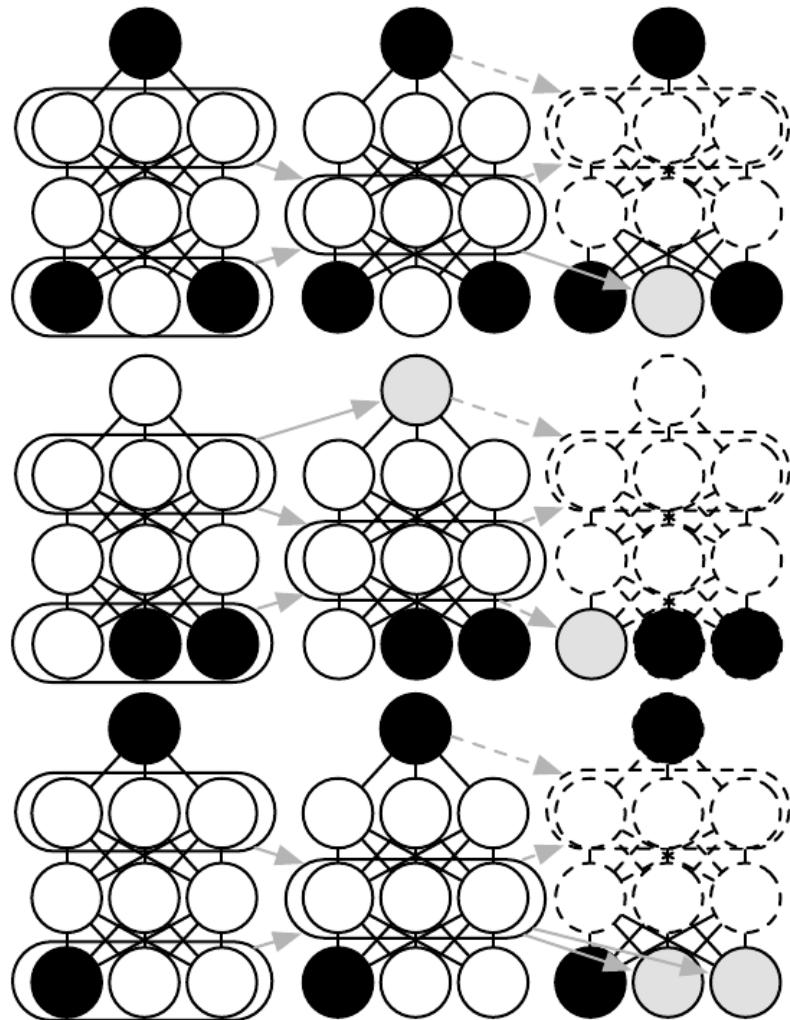


Figure 20.5: An illustration of the multi-prediction training process for a deep Boltzmann machine. Each row indicates a different example within a minibatch for the same training step. Each column represents a time step within the mean field inference process. For each example, we sample a subset of the data variables to serve as inputs to the inference process. These variables are shaded black to indicate conditioning. We then run the mean field inference process, with arrows indicating which variables influence which other variables in the process. In practical applications, we unroll mean field for several steps. In this illustration, we unroll for only two steps. Dashed arrows indicate how the process could be unrolled for more steps. The data variables that were not used as inputs to the inference process become targets, shaded in gray. We can view the inference process for each example as a recurrent network. We use gradient descent and back-propagation to train these recurrent networks to produce the correct targets given their inputs. This trains the mean field process for the MP-DBM to produce accurate estimates. Figure adapted from [Goodfellow et al. \(2013b\)](#).

The disadvantage of back-propagating through the approximate inference graph is that it does not provide a way to optimize the log-likelihood, but rather a heuristic approximation of the generalized pseudolikelihood.

The MP-DBM inspired the NADE- $k$  (Raiko *et al.*, 2014) extension to the NADE framework, which is described in Sec. 20.10.10.

The MP-DBM has some connections to dropout. Dropout shares the same parameters among many different computational graphs, with the difference between each graph being whether it includes or excludes each unit. The MP-DBM also shares parameters across many computational graphs. In the case of the MP-DBM, the difference between the graphs is whether each input unit is observed or not. When a unit is not observed, the MP-DBM does not delete it entirely as in the case of dropout. Instead, the MP-DBM treats it as a latent variable to be inferred. One could imagine applying dropout to the MP-DBM by additionally removing some units rather than making them latent.

## 20.5 Boltzmann Machines for Real-Valued Data

While Boltzmann machines were originally developed for use with binary data, many applications such as image and audio modeling seem to require the ability to represent probability distributions over real values. In some cases, it is possible to treat real-valued data in the interval  $[0, 1]$  as representing the expectation of a binary variable. For example, Hinton (2000) treats grayscale images in the training set as defining  $[0, 1]$  probability values. Each pixel defines the probability of a binary value being 1, and the binary pixels are all sampled independently from each other. This is a common procedure for evaluating binary models on grayscale image datasets. However, it is not a particularly theoretically satisfying approach, and binary images sampled independently in this way have a noisy appearance. In this section, we present Boltzmann machines that define a probability density over real-valued data.

### 20.5.1 Gaussian-Bernoulli RBMs

Restricted Boltzmann machines may be developed for many exponential family conditional distributions (Welling *et al.*, 2005). Of these, the most common is the RBM with binary hidden units and real-valued visible units, with the conditional distribution over the visible units being a Gaussian distribution whose mean is a function of the hidden units.

There are many ways of parametrizing Gaussian-Bernoulli RBMs. First, we may

choose whether to use a covariance matrix or a precision matrix for the Gaussian distribution. Here we present the precision formulation. The modification to obtain the covariance formulation is straightforward. We wish to have the conditional distribution

$$p(\mathbf{v} \mid \mathbf{h}) = \mathcal{N}(\mathbf{v}; \mathbf{W}\mathbf{h}, \boldsymbol{\beta}^{-1}). \quad (20.38)$$

We can find the terms we need to add to the energy function by expanding the unnormalized log conditional distribution:

$$\log \mathcal{N}(\mathbf{v}; \mathbf{W}\mathbf{h}, \boldsymbol{\beta}^{-1}) = -\frac{1}{2} (\mathbf{v} - \mathbf{W}\mathbf{h})^\top \boldsymbol{\beta} (\mathbf{v} - \mathbf{W}\mathbf{h}) + f(\boldsymbol{\beta}). \quad (20.39)$$

Here  $f$  encapsulates all the terms that are a function only of the parameters and not the random variables in the model. We can discard  $f$  because its only role is to normalize the distribution, and the partition function of whatever energy function we choose will carry out that role.

If we include all of the terms (with their sign flipped) involving  $\mathbf{v}$  from Eq. 20.39 in our energy function and do not add any other terms involving  $\mathbf{v}$ , then our energy function will represent the desired conditional  $p(\mathbf{v} \mid \mathbf{h})$ .

We have some freedom regarding the other conditional distribution,  $p(\mathbf{h} \mid \mathbf{v})$ . Note that Eq. 20.39 contains a term

$$\frac{1}{2} \mathbf{h}^\top \mathbf{W}^\top \boldsymbol{\beta} \mathbf{W} \mathbf{h}. \quad (20.40)$$

This term cannot be included in its entirety because it includes  $h_i h_j$  terms. These correspond to edges between the hidden units. If we included these terms, we would have a linear factor model instead of a restricted Boltzmann machine. When designing our Boltzmann machine, we simply omit these  $h_i h_j$  cross terms. Omitting them does not change the conditional  $p(\mathbf{v} \mid \mathbf{h})$  so Eq. 20.39 is still respected. However, we still have a choice about whether to include the terms involving only a single  $h_i$ . If we assume a diagonal precision matrix, we find that for each hidden unit  $h_i$  we have a term

$$\frac{1}{2} h_i \sum_j \beta_j W_{j,i}^2. \quad (20.41)$$

In the above, we used the fact that  $h_i^2 = h_i$  because  $h_i \in \{0, 1\}$ . If we include this term (with its sign flipped) in the energy function, then it will naturally bias  $h_i$  to be turned off when the weights for that unit are large and connected to visible units with high precision. The choice of whether or not to include this bias term does not affect the family of distributions the model can represent (assuming that

we include bias parameters for the hidden units) but it does affect the learning dynamics of the model. Including the term may help the hidden unit activations remain reasonable even when the weights rapidly increase in magnitude.

One way to define the energy function on a Gaussian-Bernoulli RBM is thus

$$E(\mathbf{v}, \mathbf{h}) = \frac{1}{2} \mathbf{v}^\top (\boldsymbol{\beta} \odot \mathbf{v}) - (\mathbf{v} \odot \boldsymbol{\beta})^\top \mathbf{W} \mathbf{h} - \mathbf{b}^\top \mathbf{h} \quad (20.42)$$

but we may also add extra terms or parametrize the energy in terms of the variance rather than precision if we choose.

In this derivation, we have not included a bias term on the visible units, but one could easily be added. One final source of variability in the parametrization of a Gaussian-Bernoulli RBM is the choice of how to treat the precision matrix. It may either be fixed to a constant (perhaps estimated based on the marginal precision of the data) or learned. It may also be a scalar times the identity matrix, or it may be a diagonal matrix. Typically we do not allow the precision matrix to be non-diagonal in this context, because some operations would then require inverting the matrix. In the sections ahead, we will see that other forms of Boltzmann machines permit modeling the covariance structure, using various techniques to avoid inverting the precision matrix.

### 20.5.2 Undirected Models of Conditional Covariance

While the Gaussian RBM has been the canonical energy model for real-valued data, [Ranzato \*et al.\* \(2010a\)](#) argue that the Gaussian RBM inductive bias is not well suited to the statistical variations present in some types of real-valued data, especially natural images. The problem is that much of the information content present in natural images is embedded in the covariance between pixels rather than in the raw pixel values. In other words, it is the relationships between pixels and not their absolute values where most of the useful information in images resides. Since the Gaussian RBM only models the conditional mean of the input given the hidden units, it cannot capture conditional covariance information. In response to these criticisms, alternative models have been proposed that attempt to better account for the covariance of real-valued data. These models include the mean and covariance RBM (mcRBM<sup>1</sup>), the mean-product of *t*-distribution (mPoT) model and the spike and slab RBM (ssRBM).

---

<sup>1</sup>The term “mcRBM” is pronounced by saying the name of the letters M-C-R-B-M; the “mc” is not pronounced like the “Mc” in “McDonald’s.”

**Mean and Covariance RBM** The mcRBM uses its hidden units to independently encode the conditional mean and covariance of all observed units. The mcRBM hidden layer is divided into two groups of units: mean units and covariance units. The group that models the conditional mean is simply a Gaussian RBM. The other half is a covariance RBM (Ranzato *et al.*, 2010a), also called a cRBM, whose components model the conditional covariance structure, as described below.

Specifically, with binary mean units  $\mathbf{h}^{(m)}$  and binary covariance units  $\mathbf{h}^{(c)}$ , the mcRBM model is defined as the combination of two energy functions:

$$E_{\text{mc}}(\mathbf{x}, \mathbf{h}^{(m)}, \mathbf{h}^{(c)}) = E_{\text{m}}(\mathbf{x}, \mathbf{h}^{(m)}) + E_{\text{c}}(\mathbf{x}, \mathbf{h}^{(c)}), \quad (20.43)$$

where  $E_{\text{m}}$  is the standard Gaussian-Bernoulli RBM energy function:<sup>2</sup>

$$E_{\text{m}}(\mathbf{x}, \mathbf{h}^{(m)}) = -\frac{1}{2} \mathbf{x}^\top \mathbf{x} - \sum_j \mathbf{x}^\top \mathbf{W}_{:,j} h_j^{(m)} - \sum_j b_j^{(m)} h_j^{(m)}, \quad (20.44)$$

and  $E_{\text{c}}$  is the cRBM energy function that models the conditional covariance information:

$$E_{\text{c}}(\mathbf{x}, \mathbf{h}^{(c)}) = -\frac{1}{2} \sum_j h_j^{(c)} (\mathbf{x}^\top \mathbf{r}^{(j)})^2 - \sum_j b_j^{(c)} h_j^{(c)}. \quad (20.45)$$

The parameter  $\mathbf{r}^{(j)}$  corresponds to the covariance weight vector associated with  $h_j^{(c)}$  and  $\mathbf{b}^{(c)}$  is a vector of covariance offsets. The combined energy function defines a joint distribution:

$$p_{\text{mc}}(\mathbf{x}, \mathbf{h}^{(m)}, \mathbf{h}^{(c)}) = \frac{1}{Z} \exp \left\{ -E_{\text{mc}}(\mathbf{x}, \mathbf{h}^{(m)}, \mathbf{h}^{(c)}) \right\}, \quad (20.46)$$

and a corresponding conditional distribution over the observations given  $\mathbf{h}^{(m)}$  and  $\mathbf{h}^{(c)}$  as a multivariate Gaussian distribution:

$$p_{\text{mc}}(\mathbf{x} | \mathbf{h}^{(m)}, \mathbf{h}^{(c)}) = \mathcal{N} \left( \mathbf{C}_{\mathbf{x}|\mathbf{h}}^{\text{mc}} \left( \sum_j \mathbf{W}_{:,j} h_j^{(m)} \right), \mathbf{C}_{\mathbf{x}|\mathbf{h}}^{\text{mc}} \right). \quad (20.47)$$

Note that the covariance matrix  $\mathbf{C}_{\mathbf{x}|\mathbf{h}}^{\text{mc}} = \left( \sum_j h_j \mathbf{r}^{(j)} \mathbf{r}^{(j)\top} + \mathbf{I} \right)^{-1}$  is non-diagonal and that  $\mathbf{W}$  is the weight matrix associated with the Gaussian RBM modeling the

---

<sup>2</sup>This version of the Gaussian-Bernoulli RBM energy function assumes the image data has zero mean, per pixel. Pixel offsets can easily be added to the model to account for nonzero pixel means.

conditional means. It is difficult to train the mcRBM via contrastive divergence or persistent contrastive divergence because of its non-diagonal conditional covariance structure. CD and PCD require sampling from the joint distribution of  $\mathbf{x}, \mathbf{h}^{(m)}, \mathbf{h}^{(c)}$  which, in a standard RBM, is accomplished by Gibbs sampling over the conditionals. However, in the mcRBM, sampling from  $p_{\text{mc}}(\mathbf{x} | \mathbf{h}^{(m)}, \mathbf{h}^{(c)})$  requires computing  $(\mathbf{C}^{\text{mc}})^{-1}$  at every iteration of learning. This can be an impractical computational burden for larger observations. Ranzato and Hinton (2010) avoid direct sampling from the conditional  $p_{\text{mc}}(\mathbf{x} | \mathbf{h}^{(m)}, \mathbf{h}^{(c)})$  by sampling directly from the marginal  $p(\mathbf{x})$  using Hamiltonian (hybrid) Monte Carlo (Neal, 1993) on the mcRBM free energy.

**Mean-Product of Student's  $t$ -distributions** The mean-product of Student's  $t$ -distribution (mPoT) model (Ranzato *et al.*, 2010b) extends the PoT model (Welling *et al.*, 2003a) in a manner similar to how the mcRBM extends the cRBM. This is achieved by including nonzero Gaussian means by the addition of Gaussian RBM-like hidden units. Like the mcRBM, the PoT conditional distribution over the observation is a multivariate Gaussian (with non-diagonal covariance) distribution; however, unlike the mcRBM, the complementary conditional distribution over the hidden variables is given by conditionally independent Gamma distributions. The Gamma distribution  $\mathcal{G}(k, \theta)$  is a probability distribution over positive real numbers, with mean  $k\theta$ . It is not necessary to have a more detailed understanding of the Gamma distribution to understand the basic ideas underlying the mPoT model.

The mPoT energy function is:

$$E_{\text{mPoT}}(\mathbf{x}, \mathbf{h}^{(m)}, \mathbf{h}^{(c)}) \quad (20.48)$$

$$= E_m(\mathbf{x}, \mathbf{h}^{(m)}) + \sum_j \left( h_j^{(c)} \left( 1 + \frac{1}{2} (\mathbf{r}^{(j)\top} \mathbf{x})^2 \right) + (1 - \gamma_j) \log h_j^{(c)} \right) \quad (20.49)$$

where  $\mathbf{r}^{(j)}$  is the covariance weight vector associated with unit  $h_j^{(c)}$  and  $E_m(\mathbf{x}, \mathbf{h}^{(m)})$  is as defined in Eq. 20.44.

Just as with the mcRBM, the mPoT model energy function specifies a multivariate Gaussian, with a conditional distribution over  $\mathbf{x}$  that has non-diagonal covariance. The covariance units  $\mathbf{h}^{(c)}$  are conditionally Gamma-distributed:

$$p_{\text{mPoT}}(h_j^{(c)} | \mathbf{x}) = \mathcal{G} \left( \gamma_j, 1 + \frac{1}{2} (\mathbf{v}^{(j)\top} \mathbf{x})^2 \right) \quad (20.50)$$

Learning in the mPoT model—again, like the mcRBM—is complicated by the inability to sample from the non-diagonal Gaussian conditional  $p_{\text{mPoT}}(\mathbf{x} | \mathbf{h}^{(m)}, \mathbf{h}^{(c)})$ ,

so Ranzato *et al.* (2010b) also advocate direct sampling of  $p(\mathbf{x})$  via Hamiltonian (hybrid) Monte Carlo.

**Spike and Slab Restricted Boltzmann Machines** Spike and slab restricted Boltzmann machines (Courville *et al.*, 2011) or ssRBMs provide another means of modeling the covariance structure of real-valued data. Compared to mcRBMs, ssRBMs have the advantage of requiring neither matrix inversion nor Hamiltonian Monte Carlo methods. As a model of natural images, the ssRBM is interesting in that, like the mcRBM and the mPoT model, its binary hidden units encode the conditional covariance across pixels through the use of auxiliary real-valued variables.

The spike and slab RBM has two sets of hidden units: binary *spike* units  $\mathbf{h}$ , and real-valued *slab* units  $\mathbf{s}$ . The mean of the visible units conditioned on the hidden units is given by  $(\mathbf{h} \odot \mathbf{s})\mathbf{W}^\top$ . In other words, each column  $\mathbf{W}_{:,i}$  defines a component that can appear in the input when  $h_i = 1$ . The corresponding spike variable  $h_i$  determines whether that component is present at all. The corresponding slab variable  $s_i$  determines the intensity of that component, if it is present. When a spike variable is active, the corresponding slab variable adds variance to the input along the axis defined by  $\mathbf{W}_{:,i}$ . This allows us to model the covariance of the inputs. Fortunately, contrastive divergence and persistent contrastive divergence with Gibbs sampling are still applicable. There is no need to invert any matrix.

Formally, the ssRBM model is defined via its energy function:

$$E_{ss}(\mathbf{x}, \mathbf{s}, \mathbf{h}) = - \sum_i \mathbf{x}^\top \mathbf{W}_{:,i} s_i h_i + \frac{1}{2} \mathbf{x}^\top \mathbf{\Lambda} + \sum_i \mathbf{\Phi}_i h_i \quad (20.51)$$

$$+ \frac{1}{2} \sum_i \alpha_i s_i^2 - \sum_i \alpha_i \mu_i s_i h_i - \sum_i b_i h_i + \sum_i \alpha_i \mu_i^2 h_i, \quad (20.52)$$

where  $b_i$  is the offset of the spike  $h_i$  and  $\mathbf{\Lambda}$  is a diagonal precision matrix on the observations  $\mathbf{x}$ . The parameter  $\alpha_i > 0$  is a scalar precision parameter for the real-valued slab variable  $s_i$ . The parameter  $\mathbf{\Phi}_i$  is a non-negative diagonal matrix that defines an  $\mathbf{h}$ -modulated quadratic penalty on  $\mathbf{x}$ . Each  $\mu_i$  is a mean parameter for the slab variable  $s_i$ .

With the joint distribution defined via the energy function, it is relatively straightforward to derive the ssRBM conditional distributions. For example, by marginalizing out the slab variables  $\mathbf{s}$ , the conditional distribution over the observations given the binary spike variables  $\mathbf{h}$  is given by:

$$p_{ss}(\mathbf{x} | \mathbf{h}) = \frac{1}{P(\mathbf{h})Z} \int \exp \{-E(\mathbf{x}, \mathbf{s}, \mathbf{h})\} d\mathbf{s} \quad (20.53)$$

$$= \mathcal{N} \left( \mathbf{C}_{\mathbf{x}|\mathbf{h}}^{\text{ss}} \sum_i \mathbf{W}_{:,i} \mu_i h_i, \mathbf{C}_{\mathbf{x}|\mathbf{h}}^{\text{ss}} \right) \quad (20.54)$$

where  $\mathbf{C}_{\mathbf{x}|\mathbf{h}}^{\text{ss}} = (\Lambda + \sum_i \Phi_i h_i - \sum_i \alpha_i^{-1} h_i \mathbf{W}_{:,i} \mathbf{W}_{:,i}^\top)^{-1}$ . The last equality holds only if the covariance matrix  $\mathbf{C}_{\mathbf{x}|\mathbf{h}}^{\text{ss}}$  is positive definite.

Gating by the spike variables means that the true marginal distribution over  $\mathbf{h} \odot \mathbf{s}$  is sparse. This is different from sparse coding, where samples from the model “almost never” (in the measure theoretic sense) contain zeros in the code, and MAP inference is required to impose sparsity.

Comparing the ssRBM to the mcRBM and the mPoT models, the ssRBM parametrizes the conditional covariance of the observation in a significantly different way. The mcRBM and mPoT both model the covariance structure of the observation as  $(\sum_j h_j^{(c)} \mathbf{r}^{(j)} \mathbf{r}^{(j)\top} + \mathbf{I})^{-1}$ , using the activation of the hidden units  $h_j > 0$  to enforce constraints on the conditional covariance in the direction  $\mathbf{r}^{(j)}$ . In contrast, the ssRBM specifies the conditional covariance of the observations using the hidden spike activations  $h_i = 1$  to pinch the precision matrix along the direction specified by the corresponding weight vector. The ssRBM conditional covariance is very similar to that given by a different model: the product of probabilistic principal components analysis (PoPPCA) (Williams and Agakov, 2002). In the overcomplete setting, sparse activations with the ssRBM parametrization permit significant variance (above the nominal variance given by  $\Lambda^{-1}$ ) only in the selected directions of the sparsely activated  $h_i$ . In the mcRBM or mPoT models, an overcomplete representation would mean that to capture variation in a particular direction in the observation space requires removing potentially all constraints with positive projection in that direction. This would suggest that these models are less well suited to the overcomplete setting.

The primary disadvantage of the spike and slab restricted Boltzmann machine is that some settings of the parameters can correspond to a covariance matrix that is not positive definite. Such a covariance matrix places more unnormalized probability on values that are farther from the mean, causing the integral over all possible outcomes to diverge. Generally this issue can be avoided with simple heuristic tricks. There is not yet any theoretically satisfying solution. Using constrained optimization to explicitly avoid the regions where the probability is undefined is difficult to do without being overly conservative and also preventing the model from accessing high-performing regions of parameter space.

Qualitatively, convolutional variants of the ssRBM produce excellent samples of natural images. Some examples are shown in Fig. 16.1.

The ssRBM allows for several extensions. Including higher-order interactions and average-pooling of the slab variables (Courville *et al.*, 2014) enables the model to learn excellent features for a classifier when labeled data is scarce. Adding a term to the energy function that prevents the partition function from becoming undefined results in a sparse coding model, spike and slab sparse coding (Goodfellow *et al.*, 2013d), also known as S3C.

## 20.6 Convolutional Boltzmann Machines

As seen in Chapter 9, extremely high dimensional inputs such as images place great strain on the computation, memory and statistical requirements of machine learning models. Replacing matrix multiplication by discrete convolution with a small kernel is the standard way of solving these problems for inputs that have translation invariant spatial or temporal structure. Desjardins and Bengio (2008) showed that this approach works well when applied to RBMs.

Deep convolutional networks usually require a pooling operation so that the spatial size of each successive layer decreases. Feedforward convolutional networks often use a pooling function such as the maximum of the elements to be pooled. It is unclear how to generalize this to the setting of energy-based models. We could introduce a binary pooling unit  $p$  over  $n$  binary detector units  $\mathbf{d}$  and enforce  $p = \max_i d_i$  by setting the energy function to be  $\infty$  whenever that constraint is violated. This does not scale well though, as it requires evaluating  $2^n$  different energy configurations to compute the normalization constant. For a small  $3 \times 3$  pooling region this requires  $2^9 = 512$  energy function evaluations per pooling unit!

Lee *et al.* (2009) developed a solution to this problem called *probabilistic max pooling* (not to be confused with “stochastic pooling,” which is a technique for implicitly constructing ensembles of convolutional feedforward networks). The strategy behind probabilistic max pooling is to constrain the detector units so at most one may be active at a time. This means there are only  $n + 1$  total states (one state for each of the  $n$  detector units being on, and an additional state corresponding to all of the detector units being off). The pooling unit is on if and only if one of the detector units is on. The state with all units off is assigned energy zero. We can think of this as describing a model with a single variable that has  $n + 1$  states, or equivalently as a model that has  $n + 1$  variables that assigns energy  $\infty$  to all but  $n + 1$  joint assignments of variables.

While efficient, probabilistic max pooling does force the detector units to be mutually exclusive, which may be a useful regularizing constraint in some contexts or a harmful limit on model capacity in other contexts. It also does not support

overlapping pooling regions. Overlapping pooling regions are usually required to obtain the best performance from feedforward convolutional networks, so this constraint probably greatly reduces the performance of convolutional Boltzmann machines.

[Lee et al. \(2009\)](#) demonstrated that probabilistic max pooling could be used to build convolutional deep Boltzmann machines.<sup>3</sup> This model is able to perform operations such as filling in missing portions of its input. While intellectually appealing, this model is challenging to make work in practice, and usually does not perform as well as a classifier as traditional convolutional networks trained with supervised learning.

Many convolutional models work equally well with inputs of many different spatial sizes. For Boltzmann machines, it is difficult to change the input size for a variety of reasons. The partition function changes as the size of the input changes. Moreover, many convolutional networks achieve size invariance by scaling up the size of their pooling regions proportional to the size of the input, but scaling Boltzmann machine pooling regions is awkward. Traditional convolutional neural networks can use a fixed number of pooling units and dynamically increase the size of their pooling regions in order to obtain a fixed-size representation of a variable-sized input. For Boltzmann machines, large pooling regions become too expensive for the naive approach. The approach of [Lee et al. \(2009\)](#) of making each of the detector units in the same pooling region mutually exclusive solves the computational problems, but still does not allow variable-size pooling regions. For example, suppose we learn a model with  $2 \times 2$  probabilistic max pooling over detector units that learn edge detectors. This enforces the constraint that only one of these edges may appear in each  $2 \times 2$  region. If we then increase the size of the input image by 50% in each direction, we would expect the number of edges to increase correspondingly. Instead, if we increase the size of the pooling regions by 50% in each direction to  $3 \times 3$ , then the mutual exclusivity constraint now specifies that each of these edges may only appear once in a  $3 \times 3$  region. As we grow a model’s input image in this way, the model generates edges with less density. Of course, these issues only arise when the model must use variable amounts of pooling in order to emit a fixed-size output vector. Models that use probabilistic max pooling may still accept variable-sized input images so long as the output of the model is a feature map that can scale in size proportional to the input image.

Pixels at the boundary of the image also pose some difficulty, which is exac-

---

<sup>3</sup>The publication describes the model as a “deep belief network” but because it can be described as a purely undirected model with tractable layer-wise mean field fixed point updates, it best fits the definition of a deep Boltzmann machine.

erbated by the fact that connections in a Boltzmann machine are symmetric. If we do not implicitly zero-pad the input, then there are fewer hidden units than visible units, and the visible units at the boundary of the image are not modeled well because they lie in the receptive field of fewer hidden units. However, if we do implicitly zero-pad the input, then the hidden units at the boundary are driven by fewer input pixels, and may fail to activate when needed.

## 20.7 Boltzmann Machines for Structured or Sequential Outputs

In the structured output scenario, we wish to train a model that can map from some input  $\mathbf{x}$  to some output  $\mathbf{y}$ , and the different entries of  $\mathbf{y}$  are related to each other and must obey some constraints. For example, in the speech synthesis task,  $\mathbf{y}$  is a waveform, and the entire waveform must sound like a coherent utterance.

A natural way to represent the relationships between the entries in  $\mathbf{y}$  is to use a probability distribution  $p(\mathbf{y} \mid \mathbf{x})$ . Boltzmann machines, extended to model conditional distributions, can supply this probabilistic model.

The same tool of conditional modeling with a Boltzmann machine can be used not just for structured output tasks, but also for sequence modeling. In the latter case, rather than mapping an input  $\mathbf{x}$  to an output  $\mathbf{y}$ , the model must estimate a probability distribution over a sequence of variables,  $p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)})$ . Conditional Boltzmann machines can represent factors of the form  $p(\mathbf{x}^{(t)} \mid \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t-1)})$  in order to accomplish this task.

An important sequence modeling task for the video game and film industry is modeling sequences of joint angles of skeletons used to render 3-D characters. These sequences are often collected using motion capture systems to record the movements of actors. A probabilistic model of a character’s movement allows the generation of new, previously unseen, but realistic animations. To solve this sequence modeling task, [Taylor et al. \(2007\)](#) introduced a conditional RBM modeling  $p(\mathbf{x}^{(t)} \mid \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(t-m)})$  for small  $m$ . The model is an RBM over  $p(\mathbf{x}^{(t)})$  whose bias parameters are a linear function of the preceding  $m$  values of  $\mathbf{x}$ . When we condition on different values of  $\mathbf{x}^{(t-1)}$  and earlier variables, we get a new RBM over  $\mathbf{x}$ . The weights in the RBM over  $\mathbf{x}$  never change, but by conditioning on different past values, we can change the probability of different hidden units in the RBM being active. By activating and deactivating different subsets of hidden units, we can make large changes to the probability distribution induced on  $\mathbf{x}$ . Other variants of conditional RBM ([Mnih et al., 2011](#)) and other variants of sequence

modeling using conditional RBMs are possible (Taylor and Hinton, 2009; Sutskever *et al.*, 2009; Boulanger-Lewandowski *et al.*, 2012).

Another sequence modeling task is to model the distribution over sequences of musical notes used to compose songs. Boulanger-Lewandowski *et al.* (2012) introduced the *RNN-RBM* sequence model and applied it to this task. The RNN-RBM is a generative model of a sequence of frames  $\mathbf{x}^{(t)}$  consisting of an RNN that emits the RBM parameters for each time step. Unlike the model described above, the RNN emits all of the parameters of the RBM, including the weights. To train the model, we need to be able to back-propagate the gradient of the loss function through the RNN. The loss function is not applied directly to the RNN outputs. Instead, it is applied to the RBM. This means that we must approximately differentiate the loss with respect to the RBM parameters using contrastive divergence or a related algorithm. This approximate gradient may then be back-propagated through the RNN using the usual back-propagation through time algorithm.

## 20.8 Other Boltzmann Machines

Many other variants of Boltzmann machines are possible.

Boltzmann machines may be extended with different training criteria. We have focused on Boltzmann machines trained to approximately maximize the generative criterion  $\log p(\mathbf{v})$ . It is also possible to train discriminative RBMs that aim to maximize  $\log p(y | \mathbf{v})$  instead (Larochelle and Bengio, 2008). This approach often performs the best when using a linear combination of both the generative and the discriminative criteria. Unfortunately, RBMs do not seem to be as powerful supervised learners as MLPs, at least using existing methodology.

Most Boltzmann machines used in practice have only second-order interactions in their energy functions, meaning that their energy functions are the sum of many terms and each individual term only includes the product between two random variables. An example of such a term is  $v_i W_{i,j} h_j$ . It is also possible to train higher-order Boltzmann machines (Sejnowski, 1987) whose energy function terms involve the products between many variables. Three-way interactions between a hidden unit and two different images can model spatial transformations from one frame of video to the next (Memisevic and Hinton, 2007, 2010). Multiplication by a one-hot class variable can change the relationship between visible and hidden units depending on which class is present (Nair and Hinton, 2009). One recent example of the use of higher-order interactions is a Boltzmann machine with two groups of hidden units, with one group of hidden units that interact with both the visible

units  $\mathbf{v}$  and the class label  $y$ , and another group of hidden units that interact only with the  $\mathbf{v}$  input values (Luo *et al.*, 2011). This can be interpreted as encouraging some hidden units to learn to model the input using features that are relevant to the class but also to learn extra hidden units that explain nuisance details that are necessary for the samples of  $\mathbf{v}$  to be realistic but do not determine the class of the example. Another use of higher-order interactions is to gate some features. Sohn *et al.* (2013) introduced a Boltzmann machine with third-order interactions with binary mask variables associated with each visible unit. When these masking variables are set to zero, they remove the influence of a visible unit on the hidden units. This allows visible units that are not relevant to the classification problem to be removed from the inference pathway that estimates the class.

More generally, the Boltzmann machine framework is a rich space of models permitting many more model structures than have been explored so far. Developing a new form of Boltzmann machine requires some more care and creativity than developing a new neural network layer, because it is often difficult to find an energy function that maintains tractability of all of the different conditional distributions needed to use the Boltzmann machine, but despite this required effort the field remains open to innovation.

## 20.9 Back-Propagation through Random Operations

Traditional neural networks implement a deterministic transformation of some input variables  $\mathbf{x}$ . When developing generative models, we often wish to extend neural networks to implement stochastic transformations of  $\mathbf{x}$ . One straightforward way to do this is to augment the neural network with extra inputs  $\mathbf{z}$  that are sampled from some simple probability distribution, such as a uniform or Gaussian distribution. The neural network can then continue to perform deterministic computation internally, but the function  $f(\mathbf{x}, \mathbf{z})$  will appear stochastic to an observer who does not have access to  $\mathbf{z}$ . Provided that  $f$  is continuous and differentiable, we can then compute the gradients necessary for training using back-propagation as usual.

As an example, let us consider the operation consisting of drawing samples  $y$  from a Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ :

$$y \sim \mathcal{N}(\mu, \sigma^2). \quad (20.55)$$

Because an individual sample of  $y$  is not produced by a function, but rather by a sampling process whose output changes every time we query it, it may seem counterintuitive to take the derivatives of  $y$  with respect to the parameters of

its distribution,  $\mu$  and  $\sigma^2$ . However, we can rewrite the sampling process as transforming an underlying random value  $z \sim \mathcal{N}(z; 0, 1)$  to obtain a sample from the desired distribution:

$$y = \mu + \sigma z \quad (20.56)$$

We are now able to back-propagate through the sampling operation, by regarding it as a deterministic operation with an extra input  $z$ . Crucially, the extra input is a random variable whose distribution is not a function of any of the variables whose derivatives we want to calculate. The result tells us how an infinitesimal change in  $\mu$  or  $\sigma$  would change the output if we could repeat the sampling operation again with the same value of  $z$ .

Being able to back-propagate through this sampling operation allows us to incorporate it into a larger graph. We can build elements of the graph on top of the output of the sampling distribution. For example, we can compute the derivatives of some loss function  $J(y)$ . We can also build elements of the graph whose outputs are the inputs or the parameters of the sampling operation. For example, we could build a larger graph with  $\mu = f(\mathbf{x}; \boldsymbol{\theta})$  and  $\sigma = g(\mathbf{x}; \boldsymbol{\theta})$ . In this augmented graph, we can use back-propagation through these functions to derive  $\nabla_{\boldsymbol{\theta}} J(y)$ .

The principle used in this Gaussian sampling example is more generally applicable. We can express any probability distribution of the form  $p(\mathbf{y}; \boldsymbol{\theta})$  or  $p(\mathbf{y} | \mathbf{x}; \boldsymbol{\theta})$  as  $p(\mathbf{y} | \boldsymbol{\omega})$ , where  $\boldsymbol{\omega}$  is a variable containing both parameters  $\boldsymbol{\theta}$ , and if applicable, the inputs  $\mathbf{x}$ . Given a value  $y$  sampled from distribution  $p(\mathbf{y} | \boldsymbol{\omega})$ , where  $\boldsymbol{\omega}$  may in turn be a function of other variables, we can rewrite

$$\mathbf{y} \sim p(\mathbf{y} | \boldsymbol{\omega}) \quad (20.57)$$

as

$$\mathbf{y} = f(\mathbf{z}; \boldsymbol{\omega}), \quad (20.58)$$

where  $\mathbf{z}$  is a source of randomness. We may then compute the derivatives of  $\mathbf{y}$  with respect to  $\boldsymbol{\omega}$  using traditional tools such as the back-propagation algorithm applied to  $f$ , so long as  $f$  is continuous and differentiable almost everywhere. Crucially,  $\boldsymbol{\omega}$  must not be a function of  $\mathbf{z}$ , and  $\mathbf{z}$  must not be a function of  $\boldsymbol{\omega}$ . This technique is often called the *reparametrization trick*, *stochastic back-propagation* or *perturbation analysis*.

The requirement that  $f$  be continuous and differentiable of course requires  $\mathbf{y}$  to be continuous. If we wish to back-propagate through a sampling process that produces discrete-valued samples, it may still be possible to estimate a gradient on  $\boldsymbol{\omega}$ , using reinforcement learning algorithms such as variants of the REINFORCE algorithm (Williams, 1992), discussed in Sec. 20.9.1.

In neural network applications, we typically choose  $\mathbf{z}$  to be drawn from some simple distribution, such as a unit uniform or unit Gaussian distribution, and achieve more complex distributions by allowing the deterministic portion of the network to reshape its input.

The idea of propagating gradients or optimizing through stochastic operations dates back to the mid-twentieth century (Price, 1958; Bonnet, 1964) and was first used for machine learning in the context of reinforcement learning (Williams, 1992). More recently, it has been applied to variational approximations (Opper and Archambeau, 2009) and stochastic or generative neural networks (Bengio *et al.*, 2013b; Kingma, 2013; Kingma and Welling, 2014b,a; Rezende *et al.*, 2014; Goodfellow *et al.*, 2014c). Many networks, such as denoising autoencoders or networks regularized with dropout, are also naturally designed to take noise as an input without requiring any special reparametrization to make the noise independent from the model.

### 20.9.1 Back-Propagating through Discrete Stochastic Operations

When a model emits a discrete variable  $\mathbf{y}$ , the reparametrization trick is not applicable. Suppose that the model takes inputs  $\mathbf{x}$  and parameters  $\boldsymbol{\theta}$ , both encapsulated in the vector  $\boldsymbol{\omega}$ , and combines them with random noise  $\mathbf{z}$  to produce  $\mathbf{y}$ :

$$\mathbf{y} = f(\mathbf{z}; \boldsymbol{\omega}). \quad (20.59)$$

Because  $\mathbf{y}$  is discrete,  $f$  must be a step function. The derivatives of a step function are not useful at any point. Right at each step boundary, the derivatives are undefined, but that is a small problem. The large problem is that the derivatives are zero almost everywhere, on the regions between step boundaries. The derivatives of any cost function  $J(\mathbf{y})$  therefore do not give any information for how to update the model parameters  $\boldsymbol{\theta}$ .

The REINFORCE algorithm (REward Increment = Non-negative Factor  $\times$  Offset Reinforcement  $\times$  Characteristic Eligibility) provides a framework defining a family of simple but powerful solutions (Williams, 1992). The core idea is that even though  $J(f(\mathbf{z}; \boldsymbol{\omega}))$  is a step function with useless derivatives, the expected cost  $\mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} J(f(\mathbf{z}; \boldsymbol{\omega}))$  is often a smooth function amenable to gradient descent. Although that expectation is typically not tractable when  $\mathbf{y}$  is high-dimensional (or is the result of the composition of many discrete stochastic decisions), it can be estimated without bias using a Monte Carlo average. The stochastic estimate of the gradient can be used with SGD or other stochastic gradient-based optimization techniques.

The simplest version of REINFORCE can be derived by simply differentiating the expected cost:

$$\mathbb{E}_z[J(\mathbf{y})] = \sum_{\mathbf{y}} J(\mathbf{y}) p(\mathbf{y}) \quad (20.60)$$

$$\frac{\partial \mathbb{E}[J(\mathbf{y})]}{\partial \boldsymbol{\omega}} = \sum_{\mathbf{y}} J(\mathbf{y}) \frac{\partial p(\mathbf{y})}{\partial \boldsymbol{\omega}} \quad (20.61)$$

$$= \sum_{\mathbf{y}} J(\mathbf{y}) p(\mathbf{y}) \frac{\partial \log p(\mathbf{y})}{\partial \boldsymbol{\omega}} \quad (20.62)$$

$$\approx \frac{1}{m} \sum_{\mathbf{y}^{(i)} \sim p(\mathbf{y}), i=1}^m J(\mathbf{y}^{(i)}) \frac{\partial \log p(\mathbf{y}^{(i)})}{\partial \boldsymbol{\omega}}. \quad (20.63)$$

Eq. 20.61 relies on the assumption that  $J$  does not reference  $\boldsymbol{\omega}$  directly. It is trivial to extend the approach to relax this assumption. Eq. 20.62 exploits the derivative rule for the logarithm,  $\frac{\partial \log p(\mathbf{y})}{\partial \boldsymbol{\omega}} = \frac{1}{p(\mathbf{y})} \frac{\partial p(\mathbf{y})}{\partial \boldsymbol{\omega}}$ . Eq. 20.63 gives an unbiased Monte Carlo estimator of the gradient.

Anywhere we write  $p(\mathbf{y})$  in this section, one could equally write  $p(\mathbf{y} | \mathbf{x})$ . This is because  $p(\mathbf{y})$  is parametrized by  $\boldsymbol{\omega}$ , and  $\boldsymbol{\omega}$  contains both  $\boldsymbol{\theta}$  and  $\mathbf{x}$ , if  $\mathbf{x}$  is present.

One issue with the above simple REINFORCE estimator is that it has a very high variance, so that many samples of  $\mathbf{y}$  need to be drawn to obtain a good estimator of the gradient, or equivalently, if only one sample is drawn, SGD will converge very slowly and will require a smaller learning rate. It is possible to considerably reduce the variance of that estimator by using *variance reduction* methods (Wilson, 1984; L'Ecuyer, 1994). The idea is to modify the estimator so that its expected value remains unchanged but its variance get reduced. In the context of REINFORCE, the proposed variance reduction methods involve the computation of a *baseline* that is used to offset  $J(\mathbf{y})$ . Note that any offset  $b(\boldsymbol{\omega})$  that does not depend on  $\mathbf{y}$  would not change the expectation of the estimated gradient because

$$E_{p(\mathbf{y})} \left[ \frac{\partial \log p(\mathbf{y})}{\partial \boldsymbol{\omega}} \right] = \sum_{\mathbf{y}} p(\mathbf{y}) \frac{\partial \log p(\mathbf{y})}{\partial \boldsymbol{\omega}} \quad (20.64)$$

$$= \sum_{\mathbf{y}} \frac{\partial p(\mathbf{y})}{\partial \boldsymbol{\omega}} \quad (20.65)$$

$$= \frac{\partial}{\partial \boldsymbol{\omega}} \sum_{\mathbf{y}} p(\mathbf{y}) = \frac{\partial}{\partial \boldsymbol{\omega}} 1 = 0, \quad (20.66)$$

which means that

$$E_{p(\mathbf{y})} \left[ (J(\mathbf{y}) - b(\boldsymbol{\omega})) \frac{\partial \log p(\mathbf{y})}{\partial \boldsymbol{\omega}} \right] = E_{p(\mathbf{y})} \left[ J(\mathbf{y}) \frac{\partial \log p(\mathbf{y})}{\partial \boldsymbol{\omega}} \right] - b(\boldsymbol{\omega}) E_{p(\mathbf{y})} \left[ \frac{\partial \log p(\mathbf{y})}{\partial \boldsymbol{\omega}} \right] \quad (20.67)$$

$$= E_{p(\mathbf{y})} \left[ J(\mathbf{y}) \frac{\partial \log p(\mathbf{y})}{\partial \boldsymbol{\omega}} \right]. \quad (20.68)$$

Furthermore, we can obtain the optimal  $b(\boldsymbol{\omega})$  by computing the variance of  $(J(\mathbf{y}) - b(\boldsymbol{\omega})) \frac{\partial \log p(\mathbf{y})}{\partial \boldsymbol{\omega}}$  under  $p(\mathbf{y})$  and minimizing with respect to  $b(\boldsymbol{\omega})$ . What we find is that this optimal baseline  $b^*(\boldsymbol{\omega})_i$  is different for each element  $\omega_i$  of the vector  $\boldsymbol{\omega}$ :

$$b^*(\boldsymbol{\omega})_i = \frac{E_{p(\mathbf{y})} \left[ J(\mathbf{y}) \frac{\partial \log p(\mathbf{y})}{\partial \omega_i} \right]^2}{E_{p(\mathbf{y})} \left[ \frac{\partial \log p(\mathbf{y})}{\partial \omega_i} \right]^2}. \quad (20.69)$$

The gradient estimator with respect to  $\omega_i$  then becomes

$$(J(\mathbf{y}) - b(\boldsymbol{\omega})_i) \frac{\partial \log p(\mathbf{y})}{\partial \omega_i} \quad (20.70)$$

where  $b(\boldsymbol{\omega})_i$  estimates the above  $b^*(\boldsymbol{\omega})_i$ . The estimate  $b$  is usually obtained by adding extra outputs to the neural network and training the new outputs to estimate  $E_{p(\mathbf{y})} [J(\mathbf{y}) \frac{\partial \log p(\mathbf{y})}{\partial \omega_i}]$  and  $E_{p(\mathbf{y})} \left[ \frac{\partial \log p(\mathbf{y})}{\partial \omega_i} \right]$  for each element of  $\boldsymbol{\omega}$ . These extra outputs can be trained with the mean squared error objective, using respectively  $J(\mathbf{y}) \frac{\partial \log p(\mathbf{y})}{\partial \omega_i}$  and  $\frac{\partial \log p(\mathbf{y})}{\partial \omega_i}$  as targets when  $\mathbf{y}$  is sampled from  $p(\mathbf{y})$ , for a given  $\boldsymbol{\omega}$ . The estimate  $b$  may then be recovered by substituting these estimates into Eq. 20.69. Mnih and Gregor (2014) preferred to use a single shared output (across all elements  $i$  of  $\boldsymbol{\omega}$ ) trained with the target  $J(\mathbf{y})$ , using as baseline  $b(\boldsymbol{\omega}) \approx E_{p(\mathbf{y})}[J(\mathbf{y})]$ .

Variance reduction methods have been introduced in the reinforcement learning context (Sutton *et al.*, 2000; Weaver and Tao, 2001), generalizing previous work on the case of binary reward by Dayan (1990). See Bengio *et al.* (2013b), Mnih and Gregor (2014), Ba *et al.* (2014), Mnih *et al.* (2014), or Xu *et al.* (2015) for examples of modern uses of the REINFORCE algorithm with reduced variance in the context of deep learning. In addition to the use of an input-dependent baseline  $b(\boldsymbol{\omega})$ , Mnih and Gregor (2014) found that the scale of  $(J(\mathbf{y}) - b(\boldsymbol{\omega}))$  could be adjusted during training by dividing it by its standard deviation estimated by a moving average during training, as a kind of adaptive learning rate, to counter the effect of important variations that occur during the course of training in the magnitude of this quantity. Mnih and Gregor (2014) called this heuristic *variance normalization*.

REINFORCE-based estimators can be understood as estimating the gradient by correlating choices of  $\mathbf{y}$  with corresponding values of  $J(\mathbf{y})$ . If a good value of  $\mathbf{y}$  is unlikely under the current parametrization, it might take a long time to obtain it by chance, and get the required signal that this configuration should be reinforced.

## 20.10 Directed Generative Nets

As discussed in Chapter 16, directed graphical models make up a prominent class of graphical models. While directed graphical models have been very popular within the greater machine learning community, within the smaller deep learning community they have until roughly 2013 been overshadowed by undirected models such as the RBM.

In this section we review some of the standard directed graphical models that have traditionally been associated with the deep learning community.

We have already described deep belief networks, which are a partially directed model. We have also already described sparse coding models, which can be thought of as shallow directed generative models. They are often used as feature learners in the context of deep learning, though they tend to perform poorly at sample generation and density estimation. We now describe a variety of deep, fully directed models.

### 20.10.1 Sigmoid Belief Nets

Sigmoid belief networks (Neal, 1990) are a simple form of directed graphical model with a specific kind of conditional probability distribution. In general, we can think of a sigmoid belief network as having a vector of binary states  $\mathbf{s}$ , with each element of the state influenced by its ancestors:

$$p(s_i) = \sigma \left( \sum_{j < i} W_{j,i} s_j + b_i \right). \quad (20.71)$$

The most common structure of sigmoid belief network is one that is divided into many layers, with ancestral sampling proceeding through a series of many hidden layers and then ultimately generating the visible layer. This structure is very similar to the deep belief network, except that the units at the beginning of the sampling process are independent from each other, rather than sampled from a restricted Boltzmann machine. Such a structure is interesting for a variety of

reasons. One reason is that the structure is a universal approximator of probability distributions over the visible units, in the sense that it can approximate any probability distribution over binary variables arbitrarily well, given enough depth, even if the width of the individual layers is restricted to the dimensionality of the visible layer (Sutskever and Hinton, 2008).

While generating a sample of the visible units is very efficient in a sigmoid belief network, most other operations are not. Inference over the hidden units given the visible units is intractable. Mean field inference is also intractable because the variational lower bound involves taking expectations of cliques that encompass entire layers. This problem has remained difficult enough to restrict the popularity of directed discrete networks.

One approach for performing inference in a sigmoid belief network is to construct a different lower bound that is specialized for sigmoid belief networks (Saul *et al.*, 1996). This approach has only been applied to very small networks. Another approach is to use learned inference mechanisms as described in Sec. 19.5. The Helmholtz machine (Dayan *et al.*, 1995; Dayan and Hinton, 1996) is a sigmoid belief network combined with an inference network that predicts the parameters of the mean field distribution over the hidden units. Modern approaches (Gregor *et al.*, 2014; Mnih and Gregor, 2014) to sigmoid belief networks still use this inference network approach. These techniques remain difficult due to the discrete nature of the latent variables. One cannot simply back-propagate through the output of the inference network, but instead must use the relatively unreliable machinery for back-propagating through discrete sampling processes, described in Sec. 20.9.1. Recent approaches based on importance sampling, reweighted wake-sleep (Bornstein and Bengio, 2015) and bidirectional Helmholtz machines (Bornstein *et al.*, 2015) make it possible to quickly train sigmoid belief networks and reach state-of-the-art performance on benchmark tasks.

A special case of sigmoid belief networks is the case where there are no latent variables. Learning in this case is efficient, because there is no need to marginalize latent variables out of the likelihood. A family of models called auto-regressive networks generalize this fully visible belief network to other kinds of variables besides binary variables and other structures of conditional distributions besides log-linear relationships. Auto-regressive networks are described later, in Sec. 20.10.7.

### 20.10.2 Differentiable Generator Nets

Many generative models are based on the idea of using a differentiable *generator network*. The model transforms samples of latent variables  $\mathbf{z}$  to samples  $\mathbf{x}$  or

to distributions over samples  $\mathbf{x}$  using a differentiable function  $g(\mathbf{z}; \boldsymbol{\theta}^{(g)})$  which is typically represented by a neural network. This model class includes variational autoencoders, which pair the generator net with an inference net, generative adversarial networks, which pair the generator network with a discriminator network, and techniques that train generator networks in isolation.

Generator networks are essentially just parametrized computational procedures for generating samples, where the architecture provides the family of possible distributions to sample from and the parameters select a distribution from within that family.

As an example, the standard procedure for drawing samples from a normal distribution with mean  $\boldsymbol{\mu}$  and covariance  $\boldsymbol{\Sigma}$  is to feed samples  $\mathbf{z}$  from a normal distribution with zero mean and identity covariance into a very simple generator network. This generator network contains just one affine layer:

$$\mathbf{x} = g(\mathbf{z}) = \boldsymbol{\mu} + \mathbf{L}\mathbf{z} \quad (20.72)$$

where  $\mathbf{L}$  is given by the Cholesky decomposition of  $\boldsymbol{\Sigma}$ .

Pseudorandom number generators can also use nonlinear transformations of simple distributions. For example, *inverse transform sampling* (Devroye, 2013) draws a scalar  $z$  from  $U(0,1)$  and applies a nonlinear transformation to a scalar  $x$ . In this case  $g(z)$  is given by the inverse of the cumulative distribution function  $F(x) = \int_{-\infty}^x p(v)dv$ . If we are able to specify  $p(x)$ , integrate over  $x$ , and invert the resulting function, we can sample from  $p(x)$  without using machine learning.

To generate samples from more complicated distributions that are difficult to specify directly, difficult to integrate over, or whose resulting integrals are difficult to invert, we use a feedforward network to represent a parametric family of nonlinear functions  $g$ , and use training data to infer the parameters selecting the desired function.

We can think of  $g$  as providing a nonlinear change of variables that transforms the distribution over  $\mathbf{z}$  into the desired distribution over  $\mathbf{x}$ .

Recall from Eq. 3.47 that, for invertible, differentiable, continuous  $g$ ,

$$p_z(\mathbf{z}) = p_x(g(\mathbf{z})) \left| \det\left(\frac{\partial g}{\partial \mathbf{z}}\right) \right|. \quad (20.73)$$

This implicitly imposes a probability distribution over  $\mathbf{x}$ :

$$p_x(\mathbf{x}) = \frac{p_z(g^{-1}(\mathbf{x}))}{\left| \det\left(\frac{\partial g}{\partial \mathbf{z}}\right) \right|}. \quad (20.74)$$

Of course, this formula may be difficult to evaluate, depending on the choice of  $g$ , so we often use indirect means of learning  $g$ , rather than trying to maximize  $\log p(\mathbf{x})$  directly.

In some cases, rather than using  $g$  to provide a sample of  $\mathbf{x}$  directly, we use  $g$  to define a conditional distribution over  $\mathbf{x}$ . For example, we could use a generator net whose final layer consists of sigmoid outputs to provide the mean parameters of Bernoulli distributions:

$$p(\mathbf{x}_i = 1 \mid \mathbf{z}) = g(\mathbf{z})_i. \quad (20.75)$$

In this case, when we use  $g$  to define  $p(\mathbf{x} \mid \mathbf{z})$ , we impose a distribution over  $\mathbf{x}$  by marginalizing  $\mathbf{z}$ :

$$p(\mathbf{x}) = \mathbb{E}_{\mathbf{z}} p(\mathbf{x} \mid \mathbf{z}). \quad (20.76)$$

Both approaches define a distribution  $p_g(\mathbf{x})$  and allow us to train various criteria of  $p_g$  using the reparametrization trick of Sec. 20.9.

The two different approaches to formulating generator nets—emitting the parameters of a conditional distribution versus directly emitting samples—have complementary strengths and weaknesses. When the generator net defines a conditional distribution over  $\mathbf{x}$ , it is capable of generating discrete data as well as continuous data. When the generator net provides samples directly, it is capable of generating only continuous data (we could introduce discretization in the forward propagation, but this would lose the ability to learn the model using back-propagation). The advantage to direct sampling is that we are no longer forced to use conditional distributions whose form can be easily written down and algebraically manipulated by a human designer.

Approaches based on differentiable generator networks are motivated by the success of gradient descent applied to differentiable feedforward networks for classification. In the context of supervised learning, deep feedforward networks trained with gradient-based learning seem practically guaranteed to succeed given enough hidden units and enough training data. Can this same recipe for success transfer to generative modeling?

Generative modeling seems to be more difficult than classification or regression because the learning process requires optimizing intractable criteria. In the context of differentiable generator nets, the criteria are intractable because the data does not specify both the inputs  $\mathbf{z}$  and the outputs  $\mathbf{x}$  of the generator net. In the case of supervised learning, both the inputs  $\mathbf{x}$  and the outputs  $\mathbf{y}$  were given, and the optimization procedure needs only to learn how to produce the specified mapping. In the case of generative modeling, the learning procedure needs to determine how to arrange  $\mathbf{z}$  space in a useful way and additionally how to map from  $\mathbf{z}$  to  $\mathbf{x}$ .

Dosovitskiy *et al.* (2015) studied a simplified problem, where the correspondence between  $\mathbf{z}$  and  $\mathbf{x}$  is given. Specifically, the training data is computer-rendered imagery of chairs. The latent variables  $\mathbf{z}$  are parameters given to the rendering engine describing the choice of which chair model to use, the position of the chair, and other configuration details that affect the rendering of the image. Using this synthetically generated data, a convolutional network is able to learn to map  $\mathbf{z}$  descriptions of the content of an image to  $\mathbf{x}$  approximations of rendered images. This suggests that contemporary differentiable generator networks have sufficient model capacity to be good generative models, and that contemporary optimization algorithms have the ability to fit them. The difficulty lies in determining how to train generator networks when the value of  $\mathbf{z}$  for each  $\mathbf{x}$  is not fixed and known ahead of each time.

The following sections describe several approaches to training differentiable generator nets given only training samples of  $\mathbf{x}$ .

### 20.10.3 Variational Autoencoders

The *variational autoencoder* or *VAE* (Kingma, 2013; Rezende *et al.*, 2014) is a directed model that uses learned approximate inference and can be trained purely with gradient-based methods.

To generate a sample from the model, the VAE first draws a sample  $\mathbf{z}$  from the code distribution  $p_{\text{model}}(\mathbf{z})$ . The sample is then run through a differentiable generator network  $g(\mathbf{z})$ . Finally,  $\mathbf{x}$  is sampled from a distribution  $p_{\text{model}}(\mathbf{x}; g(\mathbf{z})) = p_{\text{model}}(\mathbf{x} | \mathbf{z})$ . However, during training, the approximate inference network (or encoder)  $q(\mathbf{z} | \mathbf{x})$  is used to obtain  $\mathbf{z}$  and  $p_{\text{model}}(\mathbf{x} | \mathbf{z})$  is then viewed as a decoder network.

The key insight behind variational autoencoders is that they may be trained by maximizing the variational lower bound  $\mathcal{L}(q)$  associated with data point  $\mathbf{x}$ :

$$\mathcal{L}(q) = \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \mathbf{x})} \log p_{\text{model}}(\mathbf{z}, \mathbf{x}) + \mathcal{H}(q(\mathbf{z} | \mathbf{x})) \quad (20.77)$$

$$= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \mathbf{x})} \log p_{\text{model}}(\mathbf{x} | \mathbf{z}) - D_{\text{KL}}(q(\mathbf{z} | \mathbf{x}) || p_{\text{model}}(\mathbf{z})) \quad (20.78)$$

$$\leq \log p_{\text{model}}(\mathbf{x}). \quad (20.79)$$

In Eq. 20.77, we recognize the first term as the joint log-likelihood of the visible and hidden variables under the approximate posterior over the latent variables (just like with EM, except that we use an approximate rather than the exact posterior). We recognize also a second term, the entropy of the approximate posterior. When  $q$  is chosen to be a Gaussian distribution, with noise added to a predicted mean value, maximizing this entropy term encourages increasing the standard deviation

of this noise. More generally, this entropy term encourages the variational posterior to place high probability mass on many  $\mathbf{z}$  values that could have generated  $\mathbf{x}$ , rather than collapsing to a single point estimate of the most likely value. In Eq. 20.78, we recognize the first term as the reconstruction log-likelihood found in other autoencoders. The second term tries to make the approximate posterior distribution  $q(\mathbf{z} | \mathbf{x})$  and the model prior  $p_{\text{model}}(\mathbf{z})$  approach each other.

Traditional approaches to variational inference and learning infer  $q$  via an optimization algorithm, typically iterated fixed point equations (Sec. 19.4). These approaches are slow and often require the ability to compute  $\mathbb{E}_{\mathbf{z} \sim q} \log p_{\text{model}}(\mathbf{z}, \mathbf{x})$  in closed form. The main idea behind the variational autoencoder is to train a parametric encoder (also sometimes called an inference network or recognition model) that produces the parameters of  $q$ . So long as  $\mathbf{z}$  is a continuous variable, we can then back-propagate through samples of  $\mathbf{z}$  drawn from  $q(\mathbf{z} | \mathbf{x}) = q(\mathbf{z}; f(\mathbf{x}; \boldsymbol{\theta}))$  in order to obtain a gradient with respect to  $\boldsymbol{\theta}$ . Learning then consists solely of maximizing  $\mathcal{L}$  with respect to the parameters of the encoder and decoder. All of the expectations in  $\mathcal{L}$  may be approximated by Monte Carlo sampling.

The variational autoencoder approach is elegant, theoretically pleasing, and simple to implement. It also obtains excellent results and is among the state of the art approaches to generative modeling. Its main drawback is that samples from variational autoencoders trained on images tend to be somewhat blurry. The causes of this phenomenon are not yet known. One possibility is that the blurriness is an intrinsic effect of maximum likelihood, which minimizes  $D_{\text{KL}}(p_{\text{data}} \| p_{\text{model}})$ . As illustrated in Fig. 3.6, this means that the model will assign high probability to points that occur in the training set, but may also assign high probability to other points. These other points may include blurry images. Part of the reason that the model would choose to put probability mass on blurry images rather than some other part of the space is that the variational autoencoders used in practice usually have a Gaussian distribution for  $p_{\text{model}}(\mathbf{x}; g(\mathbf{z}))$ . Maximizing a lower bound on the likelihood of such a distribution is similar to training a traditional autoencoder with mean squared error, in the sense that it has a tendency to ignore features of the input that occupy few pixels or that cause only a small change in the brightness of the pixels that they occupy. This issue is not specific to VAEs and is shared with generative models that optimize a log-likelihood, or equivalently,  $D_{\text{KL}}(p_{\text{data}} \| p_{\text{model}})$ , as argued by Theis *et al.* (2015) and by Huszar (2015). Another troubling issue with contemporary VAE models is that they tend to use only a small subset of the dimensions of  $\mathbf{z}$ , as if the encoder was not able to transform enough of the local directions in input space to a space where the marginal distribution matches the factorized prior.

The VAE framework is very straightforward to extend to a wide range of model architectures. This is a key advantage over Boltzmann machines, which require extremely careful model design to maintain tractability. VAEs work very well with a diverse family of differentiable operators. One particularly sophisticated VAE is the *deep recurrent attention writer* or *DRAW* model (Gregor *et al.*, 2015). DRAW uses a recurrent encoder and recurrent decoder combined with an attention mechanism. The generation process for the DRAW model consists of sequentially visiting different small image patches and drawing the values of the pixels at those points. VAEs can also be extended to generate sequences by defining variational RNNs (Chung *et al.*, 2015b) by using a recurrent encoder and decoder within the VAE framework. Generating a sample from a traditional RNN involves only non-deterministic operations at the output space. Variational RNNs also have random variability at the potentially more abstract level captured by the VAE latent variables.

The VAE framework has been extended to maximize not just the traditional variational lower bound, but instead the *importance weighted autoencoder* (Burda *et al.*, 2015) objective:

$$\mathcal{L}_k(\mathbf{x}, q) = \mathbb{E}_{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(k)} \sim q(\mathbf{z} | \mathbf{x})} \left[ \log \frac{1}{k} \sum_{i=1}^k \frac{p_{\text{model}}(\mathbf{x}, \mathbf{z}^{(i)})}{q(\mathbf{z}^{(i)} | \mathbf{x})} \right]. \quad (20.80)$$

This new objective is equivalent to the traditional lower bound  $\mathcal{L}$  when  $k = 1$ . However, it may also be interpreted as forming an estimate of the true  $\log p_{\text{model}}(\mathbf{x})$  using importance sampling of  $\mathbf{z}$  from proposal distribution  $q(\mathbf{z} | \mathbf{x})$ . The importance weighted autoencoder objective is also a lower bound on  $\log p_{\text{model}}(\mathbf{x})$  and becomes tighter as  $k$  increases.

Variational autoencoders have some interesting connections to the MP-DBM and other approaches that involve back-propagation through the approximate inference graph (Goodfellow *et al.*, 2013b; Stoyanov *et al.*, 2011; Brakel *et al.*, 2013). These previous approaches required an inference procedure such as mean field fixed point equations to provide the computational graph. The variational autoencoder is defined for arbitrary computational graphs, which makes it applicable to a wider range of probabilistic model families because there is no need to restrict the choice of models to those with tractable mean field fixed point equations. The variational autoencoder also has the advantage that it increases a bound on the log-likelihood of the model, while the criteria for the MP-DBM and related models are more heuristic and have little probabilistic interpretation beyond making the results of approximate inference accurate. One disadvantage of the variational autoencoder is that it learns an inference network for only one problem, inferring  $\mathbf{z}$  given  $\mathbf{x}$ .

The older methods are able to perform approximate inference over any subset of variables given any other subset of variables, because the mean field fixed point equations specify how to share parameters between the computational graphs for all of these different problems.

One very nice property of the variational autoencoder is that simultaneously training a parametric encoder in combination with the generator network forces the model to learn a predictable coordinate system that the encoder can capture. This makes it an excellent manifold learning algorithm. See Fig. 20.6 for examples of low-dimensional manifolds learned by the variational autoencoder. In one of the cases demonstrated in the figure, the algorithm discovered two independent factors of variation present in images of faces: angle of rotation and emotional expression.

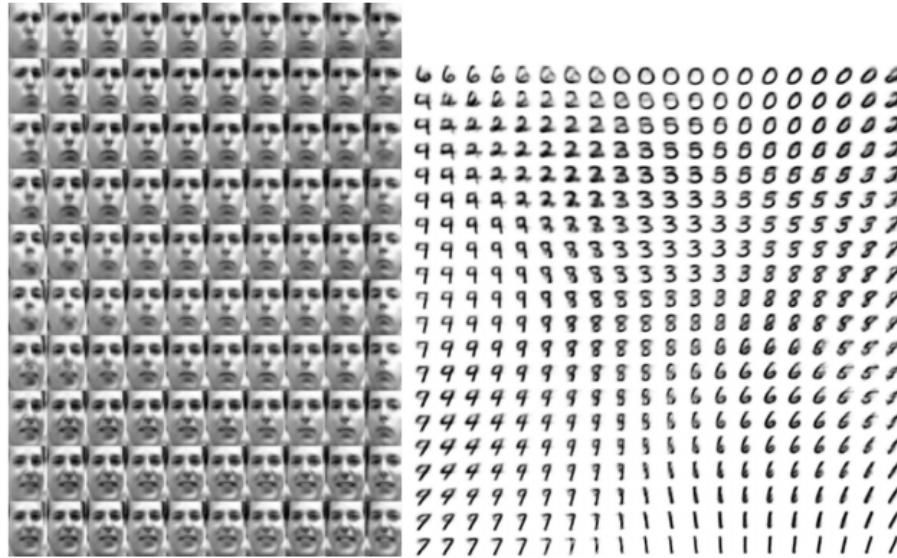


Figure 20.6: Examples of two-dimensional coordinate systems for high-dimensional manifolds, learned by a variational autoencoder (Kingma and Welling, 2014a). Two dimensions may be plotted directly on the page for visualization, so we can gain an understanding of how the model works by training a model with a 2-D latent code, even if we believe the intrinsic dimensionality of the data manifold is much higher. The images shown are not examples from the training set but images  $\mathbf{x}$  actually generated by the model  $p(\mathbf{x} | \mathbf{z})$ , simply by changing the 2-D “code”  $\mathbf{z}$  (each image corresponds to a different choice of “code”  $\mathbf{z}$  on a 2-D uniform grid). (*Left*) The two-dimensional map of the Frey faces manifold. One dimension that has been discovered (horizontal) mostly corresponds to a rotation of the face, while the other (vertical) corresponds to the emotional expression. (*Right*) The two-dimensional map of the MNIST manifold.

### 20.10.4 Generative Adversarial Networks

Generative adversarial networks or GANs (Goodfellow *et al.*, 2014c) are another generative modeling approach based on differentiable generator networks.

Generative adversarial networks are based on a game theoretic scenario in which the generator network must compete against an adversary. The generator network directly produces samples  $\mathbf{x} = g(\mathbf{z}; \boldsymbol{\theta}^{(g)})$ . Its adversary, the *discriminator network*, attempts to distinguish between samples drawn from the training data and samples drawn from the generator. The discriminator emits a probability value given by  $d(\mathbf{x}; \boldsymbol{\theta}^{(d)})$ , indicating the probability that  $\mathbf{x}$  is a real training example rather than a fake sample drawn from the model.

The simplest way to formulate learning in generative adversarial networks is as a zero-sum game, in which a function  $v(\boldsymbol{\theta}^{(g)}, \boldsymbol{\theta}^{(d)})$  determines the payoff of the discriminator. The generator receives  $-v(\boldsymbol{\theta}^{(g)}, \boldsymbol{\theta}^{(d)})$  as its own payoff. During learning, each player attempts to maximize its own payoff, so that at convergence

$$g^* = \arg \min_g \max_d v(g, d). \quad (20.81)$$

The default choice for  $v$  is

$$v(\boldsymbol{\theta}^{(g)}, \boldsymbol{\theta}^{(d)}) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log d(\mathbf{x}) + \mathbb{E}_{\mathbf{x} \sim p_{\text{model}}} \log (1 - d(\mathbf{x})). \quad (20.82)$$

This drives the discriminator to attempt to learn to correctly classify samples as real or fake. Simultaneously, the generator attempts to fool the classifier into believing its samples are real. At convergence, the generator's samples are indistinguishable from real data, and the discriminator outputs  $\frac{1}{2}$  everywhere. The discriminator may then be discarded.

The main motivation for the design of GANs is that the learning process requires neither approximate inference nor approximation of a partition function gradient. In the case where  $\max_d v(g, d)$  is convex in  $\boldsymbol{\theta}^{(g)}$  (such as the case where optimization is performed directly in the space of probability density functions) then the procedure is guaranteed to converge and is asymptotically consistent.

Unfortunately, learning in GANs can be difficult in practice when  $g$  and  $d$  are represented by neural networks and  $\max_d v(g, d)$  is not convex. Goodfellow (2014) identified non-convergence as an issue that may cause GANs to underfit. In general, simultaneous gradient descent on two players' costs is not guaranteed to reach an equilibrium. Consider for example the value function  $v(a, b) = ab$ , where one player controls  $a$  and incurs cost  $ab$ , while the other player controls  $b$  and receives a cost  $-ab$ . If we model each player as making infinitesimally small

gradient steps, each player reducing their own cost at the expense of the other player, then  $a$  and  $b$  go into a stable, circular orbit, rather than arriving at the equilibrium point at the origin. Note that the equilibria for a minimax game are not local minima of  $v$ . Instead, they are points that are simultaneously minima for both players' costs. This means that they are saddle points of  $v$  that are local minima with respect to the first player's parameters and local maxima with respect to the second player's parameters. It is possible for the two players to take turns increasing then decreasing  $v$  forever, rather than landing exactly on the saddle point where neither player is capable of reducing their cost. It is not known to what extent this non-convergence problem affects GANs.

Goodfellow (2014) identified an alternative formulation of the payoffs, in which the game is no longer zero-sum, that has the same expected gradient as maximum likelihood learning whenever the discriminator is optimal. Because maximum likelihood training converges, this reformulation of the GAN game should also converge, given enough samples. Unfortunately, this alternative formulation does not seem to perform well in practice, possibly due to suboptimality of the discriminator, or possibly due to high variance around the expected gradient.

In practice, the best-performing formulation of the GAN game is a different formulation that is neither zero-sum nor equivalent to maximum likelihood, introduced by Goodfellow *et al.* (2014c) with a heuristic motivation. In this best-performing formulation, the generator aims to increase the log probability that the discriminator makes a mistake, rather than aiming to decrease the log probability that the discriminator makes the correct prediction. This reformulation is motivated solely by the observation that it causes the derivative of the generator's cost function with respect to the discriminator's logits to remain large even in the situation where the discriminator confidently rejects all generator samples.

Stabilization of GAN learning remains an open problem. Fortunately, GAN learning performs well when the model architecture and hyperparameters are carefully selected. Radford *et al.* (2015) crafted a deep convolutional GAN (DCGAN) that performs very well for image synthesis tasks, and showed that its latent representation space captures important factors of variation, as shown in Fig. 15.9. See Fig. 20.7 for examples of images generated by a DCGAN generator.

The GAN learning problem can also be simplified by breaking the generation process into many levels of detail. It is possible to train conditional GANs (Mirza and Osindero, 2014) that learn to sample from a distribution  $p(\mathbf{x} \mid \mathbf{y})$  rather than simply sampling from a marginal distribution  $p(\mathbf{x})$ . Denton *et al.* (2015) showed that a series of conditional GANs can be trained to first generate a very low-resolution version of an image, then incrementally add details to the image.



Figure 20.7: Images generated by GANs trained on the LSUN dataset. (*Left*) Images of bedrooms generated by a DCGAN model, reproduced with permission from Radford *et al.* (2015). (*Right*) Images of churches generated by a LAPGAN model, reproduced with permission from Denton *et al.* (2015).

This technique is called the LAPGAN model, due to the use of a Laplacian pyramid to generate the images containing varying levels of detail. LAPGAN generators are able to fool not only discriminator networks but also human observers, with experimental subjects identifying up to 40% of the outputs of the network as being real data. See Fig. 20.7 for examples of images generated by a LAPGAN generator.

One unusual capability of the GAN training procedure is that it can fit probability distributions that assign zero probability to the training points. Rather than maximizing the log probability of specific points, the generator net learns to trace out a manifold whose points resemble training points in some way. Somewhat paradoxically, this means that the model may assign a log-likelihood of negative infinity to the test set, while still representing a manifold that a human observer judges to capture the essence of the generation task. This is not clearly an advantage or a disadvantage, and one may also guarantee that the generator network assigns non-zero probability to all points simply by making the last layer of the generator network add Gaussian noise to all of the generated values. Generator networks that add Gaussian noise in this manner sample from the same distribution that one obtains by using the generator network to parametrize the mean of a conditional Gaussian distribution.

Dropout seems to be important in the discriminator network. In particular, units should be stochastically dropped while computing the gradient for the generator network to follow. Following the gradient of the deterministic version of the discriminator with its weights divided by two does not seem to be as effective.

Likewise, never using dropout seems to yield poor results.

While the GAN framework is designed for differentiable generator networks, similar principles can be used to train other kinds of models. For example, *self-supervised boosting* can be used to train an RBM generator to fool a logistic regression discriminator (Welling *et al.*, 2002).

### 20.10.5 Generative Moment Matching Networks

*Generative moment matching networks* (Li *et al.*, 2015; Dziugaite *et al.*, 2015) are another form of generative model based on differentiable generator networks. Unlike VAEs and GANs, they do not need to pair the generator network with any other network—neither an inference network as used with VAEs nor a discriminator network as used with GANs.

These networks are trained with a technique called *moment matching*. The basic idea behind moment matching is to train the generator in such a way that many of the statistics of samples generated by the model are as similar as possible to those of the statistics of the examples in the training set. In this context, a *moment* is an expectation of different powers of a random variable. For example, the first moment is the mean, the second moment is the mean of the squared values, and so on. In multiple dimensions, each element of the random vector may be raised to different powers, so that a moment may be any quantity of the form

$$\mathbb{E}_{\mathbf{x}} \Pi_i x_i^{n_i} \quad (20.83)$$

where  $\mathbf{n} = [n_1, n_2, \dots, n_d]^\top$  is a vector of non-negative integers.

Upon first examination, this approach seems to be computationally infeasible. For example, if we want to match all the moments of the form  $x_i x_j$ , then we need to minimize the difference between a number of values that is quadratic in the dimension of  $\mathbf{x}$ . Moreover, even matching all of the first and second moments would only be sufficient to fit a multivariate Gaussian distribution, which captures only linear relationships between values. Our ambitions for neural networks are to capture complex nonlinear relationships, which would require far more moments. GANs avoid this problem of exhaustively enumerating all moments by using a dynamically updated discriminator, that automatically focuses its attention on whichever statistic the generator network is matching the least effectively.

Instead, generative moment matching networks can be trained by minimizing a cost function called *maximum mean discrepancy* (Schölkopf and Smola, 2002; Gretton *et al.*, 2012) or MMD. This cost function measures the error in the first moments in an infinite-dimensional space, using an implicit mapping to feature

space defined by a kernel function in order to make computations on infinite-dimensional vectors tractable. The MMD cost is zero if and only if the two distributions being compared are equal.

Visually, the samples from generative moment matching networks are somewhat disappointing. Fortunately, they can be improved by combining the generator network with an autoencoder. First, an autoencoder is trained to reconstruct the training set. Next, the encoder of the autoencoder is used to transform the entire training set into code space. The generator network is then trained to generate code samples, which may be mapped to visually pleasing samples via the decoder.

Unlike GANs, the cost function is defined only with respect to a batch of examples from both the training set and the generator network. It is not possible to make a training update as a function of only one training example or only one sample from the generator network. This is because the moments must be computed as an empirical average across many samples. When the batch size is too small, MMD can underestimate the true amount of variation in the distributions being sampled. No finite batch size is sufficiently large to eliminate this problem entirely, but larger batches reduce the amount of underestimation. When the batch size is too large, the training procedure becomes infeasibly slow, because many examples must be processed in order to compute a single small gradient step.

As with GANs, it is possible to train a generator net using MMD even if that generator net assigns zero probability to the training points.

### 20.10.6 Convolutional Generative Networks

When generating images, it is often useful to use a generator network that includes a convolutional structure (see for example Goodfellow *et al.* (2014c) or Dosovitskiy *et al.* (2015)). To do so, we use the “transpose” of the convolution operator, described in Sec. 9.5. This approach often yields more realistic images and does so using fewer parameters than using fully connected layers without parameter sharing.

Convolutional networks for recognition tasks have information flow from the image to some summarization layer at the top of the network, often a class label. As this image flows upward through the network, information is discarded as the representation of the image becomes more invariant to nuisance transformations. In a generator network, the opposite is true. Rich details must be added as the representation of the image to be generated propagates through the network, culminating in the final representation of the image, which is of course the image itself, in all of its detailed glory, with object positions and poses and textures and

lighting. The primary mechanism for discarding information in a convolutional recognition network is the pooling layer. The generator network seems to need to add information. We cannot put the inverse of a pooling layer into the generator network because most pooling functions are not invertible. A simpler operation is to merely increase the spatial size of the representation. An approach that seems to perform acceptably is to use an “un-pooling” as introduced by Dosovitskiy *et al.* (2015). This layer corresponds to the inverse of the max-pooling operation under certain simplifying conditions. First, the stride of the max-pooling operation is constrained to be equal to the width of the pooling region. Second, the maximum input within each pooling region is assumed to be the input in the upper-left corner. Finally, all non-maximal inputs within each pooling region are assumed to be zero. These are very strong and unrealistic assumptions, but they do allow the max-pooling operator to be inverted. The inverse un-pooling operation allocates a tensor of zeros, then copies each value from spatial coordinate  $i$  of the input to spatial coordinate  $i \times k$  of the output. The integer value  $k$  defines the size of the pooling region. Even though the assumptions motivating the definition of the un-pooling operator are unrealistic, the subsequent layers are able to learn to compensate for its unusual output, so the samples generated by the model as a whole are visually pleasing.

### 20.10.7 Auto-Regressive Networks

Auto-regressive networks are directed probabilistic models with no latent random variables. The conditional probability distributions in these models are represented by neural networks (sometimes extremely simple neural networks such as logistic regression). The graph structure of these models is the complete graph. They decompose a joint probability over the observed variables using the chain rule of probability to obtain a product of conditionals of the form  $P(x_d | x_{d-1}, \dots, x_1)$ . Such models have been called *fully-visible Bayes networks* (FVBMs) and used successfully in many forms, first with logistic regression for each conditional distribution (Frey, 1998) and then with neural networks with hidden units (Bengio and Bengio, 2000b; Larochelle and Murray, 2011). In some forms of auto-regressive networks, such as NADE (Larochelle and Murray, 2011), described in Sec. 20.10.10 below, we can introduce a form of parameter sharing that brings both a statistical advantage (fewer unique parameters) and a computational advantage (less computation). This is one more instance of the recurring deep learning motif of *reuse of features*.

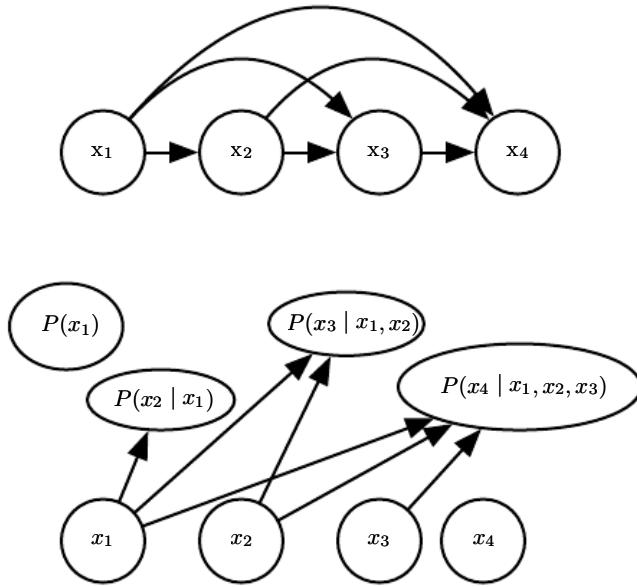


Figure 20.8: A fully visible belief network predicts the  $i$ -th variable from the  $i - 1$  previous ones. (Top) The directed graphical model for an FVBN. (Bottom) Corresponding computational graph, in the case of the logistic FVBN, where each prediction is made by a linear predictor.

### 20.10.8 Linear Auto-Regressive Networks

The simplest form of auto-regressive network has no hidden units and no sharing of parameters or features. Each  $P(x_i | x_{i-1}, \dots, x_1)$  is parametrized as a linear model (linear regression for real-valued data, logistic regression for binary data, softmax regression for discrete data). This model was introduced by Frey (1998) and has  $O(d^2)$  parameters when there are  $d$  variables to model. It is illustrated in Fig. 20.8.

If the variables are continuous, a linear auto-regressive model is merely another way to formulate a multivariate Gaussian distribution, capturing linear pairwise interactions between the observed variables.

Linear auto-regressive networks are essentially the generalization of linear classification methods to generative modeling. They therefore have the same advantages and disadvantages as linear classifiers. Like linear classifiers, they may be trained with convex loss functions, and sometimes admit closed form solutions (as in the Gaussian case). Like linear classifiers, the model itself does not offer a way of increasing its capacity, so capacity must be raised using techniques like basis expansions of the input or the kernel trick.

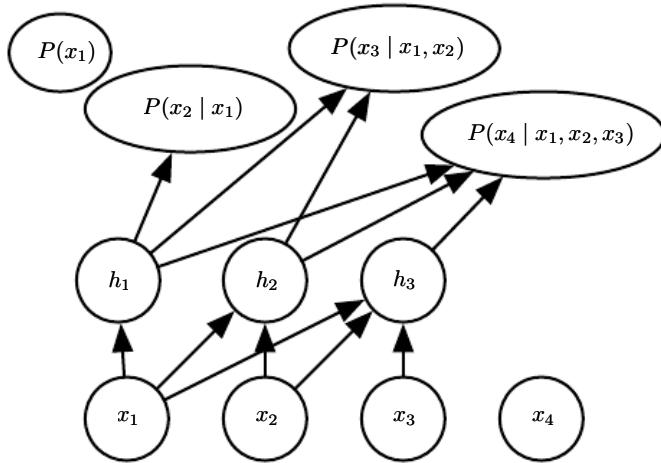


Figure 20.9: A neural auto-regressive network predicts the  $i$ -th variable  $x_i$  from the  $i - 1$  previous ones, but is parametrized so that features (groups of hidden units denoted  $h_i$ ) that are functions of  $x_1, \dots, x_i$  can be reused in predicting all of the subsequent variables  $x_{i+1}, x_{i+2}, \dots, x_d$ .

### 20.10.9 Neural Auto-Regressive Networks

Neural auto-regressive networks (Bengio and Bengio, 2000a,b) have the same left-to-right graphical model as logistic auto-regressive networks (Fig. 20.8) but employ a different parametrization of the conditional distributions within that graphical model structure. The new parametrization is more powerful in the sense that its capacity can be increased as much as needed, allowing approximation of any joint distribution. The new parametrization can also improve generalization by introducing a parameter sharing and feature sharing principle common to deep learning in general. The models were motivated by the objective of avoiding the curse of dimensionality arising out of traditional tabular graphical models, sharing the same structure as Fig. 20.8. In tabular discrete probabilistic models, each conditional distribution is represented by a table of probabilities, with one entry and one parameter for each possible configuration of the variables involved. By using a neural network instead, two advantages are obtained:

1. The parametrization of each  $P(x_i | x_{i-1}, \dots, x_1)$  by a neural network with  $(i - 1) \times k$  inputs and  $k$  outputs (if the variables are discrete and take  $k$  values, encoded one-hot) allows one to estimate the conditional probability without requiring an exponential number of parameters (and examples), yet still is able to capture high-order dependencies between the random variables.
2. Instead of having a different neural network for the prediction of each  $x_i$ ,

a *left-to-right* connectivity illustrated in Fig. 20.9 allows one to merge all the neural networks into one. Equivalently, it means that the hidden layer features computed for predicting  $x_i$  can be reused for predicting  $x_{i+k}$  ( $k > 0$ ). The hidden units are thus organized in *groups* that have the particularity that all the units in the  $i$ -th group only depend on the input values  $x_1, \dots, x_i$ . The parameters used to compute these hidden units are jointly optimized to improve the prediction of all the variables in the sequence. This is an instance of the *reuse principle* that recurs throughout deep learning in scenarios ranging from recurrent and convolutional network architectures to multi-task and transfer learning.

Each  $P(x_i | x_{i-1}, \dots, x_1)$  can represent a conditional distribution by having outputs of the neural network predict *parameters* of the conditional distribution of  $x_i$ , as discussed in Sec. 6.2.1.1. Although the original neural auto-regressive networks were initially evaluated in the context of purely discrete multivariate data (with a sigmoid output for a Bernoulli variable or softmax output for a multinoulli variable) it is natural to extend such models to continuous variables or joint distributions involving both discrete and continuous variables.

### 20.10.10 NADE

The *neural autoregressive density estimator* (NADE) is a very successful recent form of neural auto-regressive network (Larochelle and Murray, 2011). The connectivity is the same as for the original neural auto-regressive network of Bengio and Bengio (2000b) but NADE introduces an additional parameter sharing scheme, as illustrated in Fig. 20.10. The parameters of the hidden units of different groups  $j$  are shared.

The weights  $W'_{j,k,i}$  from the  $i$ -th input  $x_i$  to the  $k$ -th element of the  $j$ -th group of hidden unit  $h_k^{(j)}$  ( $j \geq i$ ) are shared among the groups:

$$W'_{j,k,i} = W_{k,i}. \quad (20.84)$$

The remaining weights, where  $j < i$ , are zero.

Larochelle and Murray (2011) chose this sharing scheme so that forward propagation in a NADE model loosely resembles the computations performed in mean field inference to fill in missing inputs in an RBM. This mean field inference corresponds to running a recurrent network with shared weights and the first step of that inference is the same as in NADE. The only difference is that with NADE, the output weights connecting the hidden units to the output are parametrized

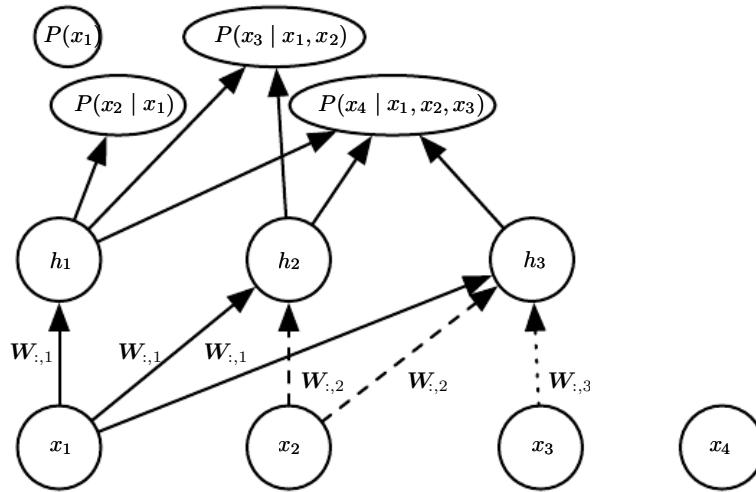


Figure 20.10: An illustration of the neural autoregressive density estimator (NADE). The hidden units are organized in groups  $\mathbf{h}^{(j)}$  so that only the inputs  $x_1, \dots, x_i$  participate in computing  $\mathbf{h}^{(i)}$  and predicting  $P(x_j | x_{j-1}, \dots, x_1)$ , for  $j > i$ . NADE is differentiated from earlier neural auto-regressive networks by the use of a particular weight sharing pattern:  $W'_{j,k,i} = W_{k,i}$  is shared (indicated in the figure by the use of the same line pattern for every instance of a replicated weight) for all the weights going out from  $x_i$  to the  $k$ -th unit of any group  $j \geq i$ . Recall that the vector  $(W_{1,i}, W_{2,i}, \dots, W_{n,i})$  is denoted  $\mathbf{W}_{:,i}$ .

independently from the weights connecting the input units to the hidden units. In the RBM, the hidden-to-output weights are the transpose of the input-to-hidden weights. The NADE architecture can be extended to mimic not just one time step of the mean field recurrent inference but to mimic  $k$  steps. This approach is called NADE- $k$  (Raiko *et al.*, 2014).

As mentioned previously, auto-regressive networks may be extend to process continuous-valued data. A particularly powerful and generic way of parametrizing a continuous density is as a Gaussian mixture (introduced in Sec. 3.9.6) with mixture weights  $\alpha_i$  (the coefficient or prior probability for component  $i$ ), per-component conditional mean  $\mu_i$  and per-component conditional variance  $\sigma_i^2$ . A model called RNADE (Uria *et al.*, 2013) uses this parametrization to extend NADE to real values. As with other mixture density networks, the parameters of this distribution are outputs of the network, with the mixture weight probabilities produced by a softmax unit, and the variances parametrized so that they are positive. Stochastic gradient descent can be numerically ill-behaved due to the interactions between the conditional means  $\mu_i$  and the conditional variances  $\sigma_i^2$ . To reduce this difficulty, Uria *et al.* (2013) use a pseudo-gradient that replaces the gradient on the mean, in the back-propagation phase.

Another very interesting extension of the neural auto-regressive architectures gets rid of the need to choose an arbitrary *order* for the observed variables ([Murray and Larochelle, 2014](#)). In auto-regressive networks, the idea is to train the network to be able to cope with any order by randomly sampling orders and providing the information to hidden units specifying which of the inputs are observed (on the right side of the conditioning bar) and which are to be predicted and are thus considered missing (on the left side of the conditioning bar). This is nice because it allows one to use a trained auto-regressive network to *perform any inference problem* (i.e. predict or sample from the probability distribution over any subset of variables given any subset) extremely efficiently. Finally, since many orders of variables are possible ( $n!$  for  $n$  variables) and each order  $o$  of variables yields a different  $p(\mathbf{x} | o)$ , we can form an ensemble of models for many values of  $o$ :

$$p_{\text{ensemble}}(\mathbf{x}) = \frac{1}{k} \sum_{i=1}^k p(\mathbf{x} | o^{(i)}). \quad (20.85)$$

This ensemble model usually generalizes better and assigns higher probability to the test set than does an individual model defined by a single ordering.

In the same paper, the authors propose deep versions of the architecture, but unfortunately that immediately makes computation as expensive as in the original neural auto-regressive neural network ([Bengio and Bengio, 2000b](#)). The first layer and the output layer can still be computed in  $O(nh)$  multiply-add operations, as in the regular NADE, where  $h$  is the number of hidden units (the size of the groups  $h_i$ , in Fig. 20.10 and Fig. 20.9), whereas it is  $O(n^2h)$  in [Bengio and Bengio \(2000b\)](#). However, for the other hidden layers, the computation is  $O(n^2h^2)$  if every “previous” group at layer  $l$  participates in predicting the “next” group at layer  $l+1$ , assuming  $n$  groups of  $h$  hidden units at each layer. Making the  $i$ -th group at layer  $l+1$  only depend on the  $i$ -th group, as in [Murray and Larochelle \(2014\)](#) at layer  $l$  reduces it to  $O(nh^2)$ , which is still  $h$  times worse than the regular NADE.

## 20.11 Drawing Samples from Autoencoders

In Chapter 14, we saw that many kinds of autoencoders learn the data distribution. There are close connections between score matching, denoising autoencoders, and contractive autoencoders. These connections demonstrate that some kinds of autoencoders learn the data distribution in some way. We have not yet seen how to draw samples from such models.

Some kinds of autoencoders, such as the variational autoencoder, explicitly

represent a probability distribution and admit straightforward ancestral sampling. Most other kinds of autoencoders require MCMC sampling.

Contractive autoencoders are designed to recover an estimate of the tangent plane of the data manifold. This means that repeated encoding and decoding with injected noise will induce a random walk along the surface of the manifold (Rifai *et al.*, 2012; Mesnil *et al.*, 2012). This manifold diffusion technique is a kind of Markov chain.

There is also a more general Markov chain that can sample from any denoising autoencoder.

### 20.11.1 Markov Chain Associated with any Denoising Autoencoder

The above discussion left open the question of what noise to inject and where, in order to obtain a Markov chain that would generate from the distribution estimated by the autoencoder. Bengio *et al.* (2013c) showed how to construct such a Markov chain for *generalized denoising autoencoders*. Generalized denoising autoencoders are specified by a denoising distribution for sampling an estimate of the clean input given the corrupted input.

Each step of the Markov chain that generates from the estimated distribution consists of the following sub-steps, illustrated in Fig. 20.11:

1. Starting from the previous state  $\mathbf{x}$ , inject corruption noise, sampling  $\tilde{\mathbf{x}}$  from  $C(\tilde{\mathbf{x}} | \mathbf{x})$ .
2. Encode  $\tilde{\mathbf{x}}$  into  $\mathbf{h} = f(\tilde{\mathbf{x}})$ .
3. Decode  $\mathbf{h}$  to obtain the parameters  $\omega = g(\mathbf{h})$  of  $p(\mathbf{x} | \omega = g(\mathbf{h})) = p(\mathbf{x} | \tilde{\mathbf{x}})$ .
4. Sample the next state  $\mathbf{x}$  from  $p(\mathbf{x} | \omega = g(\mathbf{h})) = p(\mathbf{x} | \tilde{\mathbf{x}})$ .

Bengio *et al.* (2014) showed that if the autoencoder  $p(\mathbf{x} | \tilde{\mathbf{x}})$  forms a consistent estimator of the corresponding true conditional distribution, then the stationary distribution of the above Markov chain forms a consistent estimator (albeit an implicit one) of the data generating distribution of  $\mathbf{x}$ .

### 20.11.2 Clamping and Conditional Sampling

Similarly to Boltzmann machines, denoising autoencoders and their generalizations (such as GSNs, described below) can be used to sample from a conditional distribution  $p(\mathbf{x}_f | \mathbf{x}_o)$ , simply by clamping the *observed* units  $\mathbf{x}_f$  and only resampling

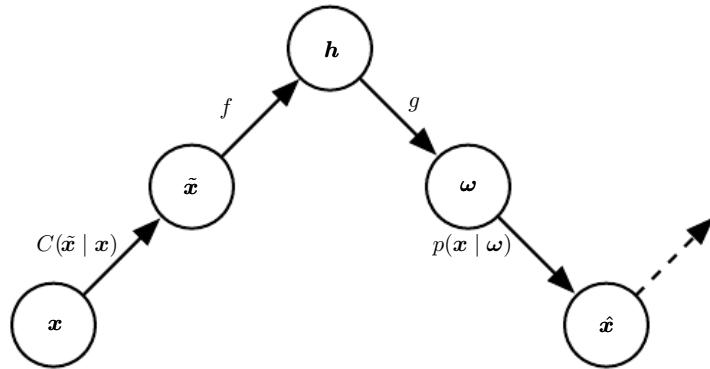


Figure 20.11: Each step of the Markov chain associated with a trained denoising autoencoder, that generates the samples from the probabilistic model implicitly trained by the denoising log-likelihood criterion. Each step consists in (a) injecting noise via corruption process  $C$  in state  $\mathbf{x}$ , yielding  $\tilde{\mathbf{x}}$ , (b) encoding it with function  $f$ , yielding  $\mathbf{h} = f(\tilde{\mathbf{x}})$ , (c) decoding the result with function  $g$ , yielding parameters  $\omega$  for the reconstruction distribution  $p(\mathbf{x} | \omega)$ , and (d) given  $\omega$ , sampling a new state from the reconstruction distribution  $p(\mathbf{x} | \omega = g(f(\tilde{\mathbf{x}})))$ . In the typical squared reconstruction error case,  $g(\mathbf{h}) = \hat{\mathbf{x}}$ , which estimates  $\mathbb{E}[\mathbf{x} | \tilde{\mathbf{x}}]$ , corruption consists in adding Gaussian noise and sampling from  $p(\mathbf{x} | \omega)$  consists in adding Gaussian noise, a second time, to the reconstruction  $\hat{\mathbf{x}}$ . The latter noise level should correspond to the mean squared error of reconstructions, whereas the injected noise is a hyperparameter that controls the mixing speed as well as the extent to which the estimator smooths the empirical distribution (Vincent, 2011). In the example illustrated here, only the  $C$  and  $p$  conditionals are stochastic steps ( $f$  and  $g$  are deterministic computations), although noise can also be injected inside the autoencoder, as in generative stochastic networks (Bengio *et al.*, 2014).

the *free* units  $\mathbf{x}_o$  given  $\mathbf{x}_f$  and the sampled latent variables (if any). For example, MP-DBMs can be interpreted as a form of denoising autoencoder, and are able to sample missing inputs. GSNs later generalized some of the ideas present in MP-DBMs to perform the same operation (Bengio *et al.*, 2014). Alain *et al.* (2015) identified a missing condition from Proposition 1 of Bengio *et al.* (2014), which is that the transition operator (defined by the stochastic mapping going from one state of the chain to the next) should satisfy a property called *detailed balance*, which specifies that a Markov Chain at equilibrium will remain in equilibrium whether the transition operator is run in forward or reverse.

An experiment in clamping half of the pixels (the right part of the image) and running the Markov chain on the other half is shown in Fig. 20.12.



Figure 20.12: Illustration of clamping the right half of the image and running the Markov Chain by resampling only the left half at each step. These samples come from a GSN trained to reconstruct MNIST digits at each time step using the walkback procedure.

### 20.11.3 Walk-Back Training Procedure

The walk-back training procedure was proposed by Bengio *et al.* (2013c) as a way to accelerate the convergence of generative training of denoising autoencoders. Instead of performing a one-step encode-decode reconstruction, this procedure consists in alternative multiple stochastic encode-decode steps (as in the generative

Markov chain) initialized at a training example (just like with the contrastive divergence algorithm, described in Sec. 18.2) and penalizing the last probabilistic reconstructions (or all of the reconstructions along the way).

Training with  $k$  steps is equivalent (in the sense of achieving the same stationary distribution) as training with one step, but practically has the advantage that spurious modes farther from the data can be removed more efficiently.

## 20.12 Generative Stochastic Networks

*Generative stochastic networks* or *GSNs* (Bengio *et al.*, 2014) are generalizations of denoising autoencoders that include latent variables  $\mathbf{h}$  in the generative Markov chain, in addition to the visible variables (usually denoted  $\mathbf{x}$ ).

A GSN is parametrized by two conditional probability distributions which specify one step of the Markov chain:

1.  $p(\mathbf{x}^{(k)} | \mathbf{h}^{(k)})$  tells how to generate the next visible variable given the current latent state. Such a “reconstruction distribution” is also found in denoising autoencoders, RBMs, DBNs and DBMs.
2.  $p(\mathbf{h}^{(k)} | \mathbf{h}^{(k-1)}, \mathbf{x}^{(k-1)})$  tells how to update the latent state variable, given the previous latent state and visible variable.

Denoising autoencoders and GSNs differ from classical probabilistic models (directed or undirected) in that they parametrize the generative process itself rather than the mathematical specification of the joint distribution of visible and latent variables. Instead, the latter is defined **implicitly, if it exists**, as the stationary distribution of the generative Markov chain. The conditions for existence of the stationary distribution are mild and are the same conditions required by standard MCMC methods (see Sec. 17.3). These conditions are necessary to guarantee that the chain mixes, but they can be violated by some choices of the transition distributions (for example, if they were deterministic).

One could imagine different training criteria for GSNs. The one proposed and evaluated by Bengio *et al.* (2014) is simply reconstruction log-probability on the visible units, just like for denoising autoencoders. This is achieved by clamping  $\mathbf{x}^{(0)} = \mathbf{x}$  to the observed example and maximizing the probability of generating  $\mathbf{x}$  at some subsequent time steps, i.e., maximizing  $\log p(\mathbf{x}^{(k)} = \mathbf{x} | \mathbf{h}^{(k)})$ , where  $\mathbf{h}^{(k)}$  is sampled from the chain, given  $\mathbf{x}^{(0)} = \mathbf{x}$ . In order to estimate the gradient of  $\log p(\mathbf{x}^{(k)} = \mathbf{x} | \mathbf{h}^{(k)})$  with respect to the other pieces of the model, Bengio *et al.* (2014) use the reparametrization trick, introduced in Sec. 20.9.

The *walk-back training* protocol (described in Sec. 20.11.3) was used (Bengio *et al.*, 2014) to improve training convergence of GSNs.

### 20.12.1 Discriminant GSNs

The original formulation of GSNs (Bengio *et al.*, 2014) was meant for unsupervised learning and implicitly modeling  $p(\mathbf{x})$  for observed data  $\mathbf{x}$ , but it is possible to modify the framework to optimize  $p(\mathbf{y} \mid \mathbf{x})$ .

For example, Zhou and Troyanskaya (2014) generalize GSNs in this way, by only back-propagating the reconstruction log-probability over the output variables, keeping the input variables fixed. They applied this successfully to model **sequences** (protein secondary structure) and introduced a (one-dimensional) **convolutional** structure in the transition operator of the Markov chain. It is important to remember that, for each step of the Markov chain, one generates a new sequence for each layer, and that sequence is the input for computing other layer values (say the one below and the one above) at the next time step.

Hence the Markov chain is really over the **output variable** (and associated higher-level hidden layers), and the input sequence only serves to condition that chain, with back-propagation allowing to learn how the input sequence can condition the output distribution implicitly represented by the Markov chain. It is therefore a case of using the GSN in the context of **structured outputs**, where  $p(\mathbf{y} \mid \mathbf{x})$  does not have a simple parametric form but instead the components of  $\mathbf{y}$  are statistically dependent of each other, given  $\mathbf{x}$ , in complicated ways.

Zöhrer and Pernkopf (2014) introduced a hybrid model that combines a supervised objective (as in the above work) and an unsupervised objective (as in the original GSN work), by simply adding (with a different weight) the supervised and unsupervised costs i.e., the reconstruction log-probabilities of  $\mathbf{y}$  and  $\mathbf{x}$  respectively. Such a hybrid criterion had previously been introduced for RBMs by Larochelle and Bengio (2008). They show improved classification performance using this scheme.

## 20.13 Other Generation Schemes

The methods we have described so far use either MCMC sampling, ancestral sampling, or some mixture of the two to generate samples. While these are the most popular approaches to generative modeling, they are by no means the only approaches.

Sohl-Dickstein *et al.* (2015) developed a *diffusion inversion* training scheme for learning a generative model, based on non-equilibrium thermodynamics. The approach is based on the idea that the probability distributions we wish to sample from have structure. This structure can gradually be destroyed by a diffusion process that incrementally changes the probability distribution to have more entropy. To form a generative model, we can run the process in reverse, by training a model that gradually restores the structure to an unstructured distribution. By iteratively applying a process that brings a distribution closer to the target one, we can gradually approach that target distribution. This approach resembles MCMC methods in the sense that it involves many iterations to produce a sample. However, the model is defined to be the probability distribution produced by the final step of the chain. In this sense, there is no approximation induced by the iterative procedure. The approach introduced by Sohl-Dickstein *et al.* (2015) is also very close to the generative interpretation of the denoising autoencoder (Sec. 20.11.1). Like with the denoising autoencoder, the training objective trains a transition operator which attempts to probabilistically undo the effect of adding some noise, trying to undo one step of the diffusion process. If we compare with the walkback training procedure (Sec. 20.11.3) for denoising autoencoders and GSNs, the main difference is that instead of reconstructing only towards the observed training point  $\mathbf{x}$ , the objective function only tries to reconstruct towards the previous point in the diffusion trajectory that started at  $\mathbf{x}$  (which should be easier). This addresses the following dilemma present with the ordinary reconstruction log-likelihood objective of denoising autoencoders: with small levels of noise the learner only sees configurations near the data points, while with large levels of noise it is asked to do an almost impossible job (because the denoising distribution is going to be highly complex and multi-modal). With the diffusion inversion objective, the learner can learn more precisely the shape of the density around the data points as well as remove spurious modes that could show up far from the data points.

Another approach to sample generation is the *approximate Bayesian computation* (ABC) framework (Rubin *et al.*, 1984). In this approach, samples are rejected or modified in order to make the moments of selected functions of the samples match those of the desired distribution. While this idea uses the moments of the samples like in moment matching, it is different from moment matching because it modifies the samples themselves, rather than training the model to automatically emit samples with the correct moments. Bachman and Precup (2015) showed how to use ideas from ABC in the context of deep learning, by using ABC to shape the MCMC trajectories of GSNs.

We expect that many other possible approaches to generative modeling await discovery.

## 20.14 Evaluating Generative Models

Researchers studying generative models often need to compare one generative model to another, usually in order to demonstrate that a newly invented generative model is better at capturing some distribution than the pre-existing models.

This can be a difficult and subtle task. In many cases, we can not actually evaluate the log probability of the data under the model, but only an approximation. In these cases, it is important to think and communicate clearly about exactly what is being measured. For example, suppose we can evaluate a stochastic estimate of the log-likelihood for model A, and a deterministic lower bound on the log-likelihood for model B. If model A gets a higher score than model B, which is better? If we care about determining which model has a better internal representation of the distribution, we actually cannot tell, unless we have some way of determining how loose the bound for model B is. However, if we care about how well we can use the model in practice, for example to perform anomaly detection, then it is fair to say that a model is preferable based on a criterion specific to the practical task of interest, e.g., based on ranking test examples and ranking criteria such as precision and recall.

Another subtlety of evaluating generative models is that the evaluation metrics are often hard research problems in and of themselves. It can be very difficult to establish that models are being compared fairly. For example, suppose we use AIS to estimate  $\log Z$  in order to compute  $\log \tilde{p}(\mathbf{x}) - \log Z$  for a new model we have just invented. A computationally economical implementation of AIS may fail to find several modes of the model distribution and underestimate  $Z$ , which will result in us overestimating  $\log p(\mathbf{x})$ . It can thus be difficult to tell whether a high likelihood estimate is due to a good model or a bad AIS implementation.

Other fields of machine learning usually allow for some variation in the preprocessing of the data. For example, when comparing the accuracy of object recognition algorithms, it is usually acceptable to preprocess the input images slightly differently for each algorithm based on what kind of input requirements it has. Generative modeling is different because changes in preprocessing, even very small and subtle ones, are completely unacceptable. Any change to the input data changes the distribution to be captured and fundamentally alters the task. For example, multiplying the input by 0.1 will artificially increase likelihood by a factor of 10.

Issues with preprocessing commonly arise when benchmarking generative models on the MNIST dataset, one of the more popular generative modeling benchmarks. MNIST consists of grayscale images. Some models treat MNIST images as points

in a real vector space, while others treat them as binary. Yet others treat the grayscale values as probabilities for a binary samples. It is essential to compare real-valued models only to other real-valued models and binary-valued models only to other binary-valued models. Otherwise the likelihoods measured are not on the same space. For binary-valued models, the log-likelihood can be at most zero, while for real-valued models it can be arbitrarily high, since it is the measurement of a density. Among binary models, it is important to compare models using exactly the same kind of binarization. For example, we might binarize a gray pixel to 0 or 1 by thresholding at 0.5, or by drawing a random sample whose probability of being 1 is given by the gray pixel intensity. If we use the random binarization, we might binarize the whole dataset once, or we might draw a different random example for each step of training and then draw multiple samples for evaluation. Each of these three schemes yields wildly different likelihood numbers, and when comparing different models it is important that both models use the same binarization scheme for training and for evaluation. In fact, researchers who apply a single random binarization step share a file containing the results of the random binarization, so that there is no difference in results based on different outcomes of the binarization step.

Because being able to generate realistic samples from the data distribution is one of the goals of a generative model, practitioners often evaluate generative models by visually inspecting the samples. In the best case, this is done not by the researchers themselves, but by experimental subjects who do not know the source of the samples (Denton *et al.*, 2015). Unfortunately, it is possible for a very poor probabilistic model to produce very good samples. A common practice to verify if the model only copies some of the training examples is illustrated in Fig. 16.1. The idea is to show for some of the generated samples their nearest neighbor in the training set, according to Euclidean distance in the space of  $\mathbf{x}$ . This test is intended to detect the case where the model overfits the training set and just reproduces training instances. It is even possible to simultaneously underfit and overfit yet still produce samples that individually look good. Imagine a generative model trained on images of dogs and cats that simply learns to reproduce the training images of dogs. Such a model has clearly overfit, because it does not produce images that were not in the training set, but it has also underfit, because it assigns no probability to the training images of cats. Yet a human observer would judge each individual image of a dog to be high quality. In this simple example, it would be easy for a human observer who can inspect many samples to determine that the cats are absent. In more realistic settings, a generative model trained on data with tens of thousands of modes may ignore a small number of modes, and a human observer would not easily be able to inspect or remember

enough images to detect the missing variation.

Since the visual quality of samples is not a reliable guide, we often also evaluate the log-likelihood that the model assigns to the test data, when this is computationally feasible. Unfortunately, in some cases the likelihood seems not to measure any attribute of the model that we really care about. For example, real-valued models of MNIST can obtain arbitrarily high likelihood by assigning arbitrarily low variance to background pixels that never change. Models and algorithms that detect these constant features can reap unlimited rewards, even though this is not a very useful thing to do. The potential to achieve a cost approaching negative infinity is present for any kind of maximum likelihood problem with real values, but it is especially problematic for generative models of MNIST because so many of the output values are trivial to predict. This strongly suggests a need for developing other ways of evaluating generative models.

Theis *et al.* (2015) review many of the issues involved in evaluating generative models, including many of the ideas described above. They highlight the fact that there are many different uses of generative models and that the choice of metric must match the intended use of the model. For example, some generative models are better at assigning high probability to most realistic points while other generative models are better at rarely assigning high probability to unrealistic points. These differences can result from whether a generative model is designed to minimize  $D_{\text{KL}}(p_{\text{data}} \| p_{\text{model}})$  or  $D_{\text{KL}}(p_{\text{model}} \| p_{\text{data}})$ , as illustrated in Fig. 3.6. Unfortunately, even when we restrict the use of each metric to the task it is most suited for, all of the metrics currently in use continue to have serious weaknesses. One of the most important research topics in generative modeling is therefore not just how to improve generative models, but in fact, designing new techniques to measure our progress.

## 20.15 Conclusion

Training generative models with hidden units is a powerful way to make models understand the world represented in the given training data. By learning a model  $p_{\text{model}}(\mathbf{x})$  and a representation  $p_{\text{model}}(\mathbf{h} \mid \mathbf{x})$ , a generative model can provide answers to many inference problems about the relationships between input variables in  $\mathbf{x}$  and can provide many different ways of representing  $\mathbf{x}$  by taking expectations of  $\mathbf{h}$  at different layers of the hierarchy. Generative models hold the promise to provide AI systems with a framework for all of the many different intuitive concepts they need to understand, and the ability to reason about these concepts in the face of uncertainty. We hope that our readers will find new ways to make these

approaches more powerful and continue the journey to understanding the principles that underlie learning and intelligence.

# Bibliography

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., and Zheng, X. (2015). TensorFlow: Large-scale machine learning on heterogeneous systems. Software available from tensorflow.org. [25, 212, 449](#)
- Ackley, D. H., Hinton, G. E., and Sejnowski, T. J. (1985). A learning algorithm for Boltzmann machines. *Cognitive Science*, **9**, 147–169. [573, 656](#)
- Alain, G. and Bengio, Y. (2013). What regularized auto-encoders learn from the data generating distribution. In *ICLR'2013, arXiv:1211.4246*. [510, 516, 524](#)
- Alain, G., Bengio, Y., Yao, L., Éric Thibodeau-Laufer, Yosinski, J., and Vincent, P. (2015). GSNs: Generative stochastic networks. *arXiv:1503.05571*. [513, 715](#)
- Anderson, E. (1935). The Irises of the Gaspé Peninsula. *Bulletin of the American Iris Society*, **59**, 2–5. [21](#)
- Ba, J., Mnih, V., and Kavukcuoglu, K. (2014). Multiple object recognition with visual attention. *arXiv:1412.7755*. [693](#)
- Bachman, P. and Precup, D. (2015). Variational generative stochastic networks with collaborative shaping. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, pages 1964–1972. [718](#)
- Bacon, P.-L., Bengio, E., Pineau, J., and Precup, D. (2015). Conditional computation in neural networks using a decision-theoretic approach. In *2nd Multidisciplinary Conference on Reinforcement Learning and Decision Making (RLDM 2015)*. [453](#)
- Bagnell, J. A. and Bradley, D. M. (2009). Differentiable sparse coding. In D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, editors, *Advances in Neural Information Processing Systems 21 (NIPS'08)*, pages 113–120. [501](#)

- Bahdanau, D., Cho, K., and Bengio, Y. (2015). Neural machine translation by jointly learning to align and translate. In *ICLR'2015*, *arXiv:1409.0473*. [25](#), [101](#), [399](#), [421](#), [422](#), [468](#), [478](#)
- Bahl, L. R., Brown, P., de Souza, P. V., and Mercer, R. L. (1987). Speech recognition with continuous-parameter hidden Markov models. *Computer, Speech and Language*, **2**, 219–234. [461](#)
- Baldi, P. and Hornik, K. (1989). Neural networks and principal component analysis: Learning from examples without local minima. *Neural Networks*, **2**, 53–58. [286](#)
- Baldi, P., Brunak, S., Frasconi, P., Soda, G., and Pollastri, G. (1999). Exploiting the past and the future in protein secondary structure prediction. *Bioinformatics*, **15**(11), 937–946. [396](#)
- Baldi, P., Sadowski, P., and Whiteson, D. (2014). Searching for exotic particles in high-energy physics with deep learning. *Nature communications*, **5**. [26](#)
- Ballard, D. H., Hinton, G. E., and Sejnowski, T. J. (1983). Parallel vision computation. *Nature*. [455](#)
- Barlow, H. B. (1989). Unsupervised learning. *Neural Computation*, **1**, 295–311. [146](#)
- Barron, A. E. (1993). Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Trans. on Information Theory*, **39**, 930–945. [198](#)
- Bartholomew, D. J. (1987). *Latent variable models and factor analysis*. Oxford University Press. [493](#)
- Basilevsky, A. (1994). *Statistical Factor Analysis and Related Methods: Theory and Applications*. Wiley. [493](#)
- Bastien, F., Lamblin, P., Pascanu, R., Bergstra, J., Goodfellow, I. J., Bergeron, A., Bouchard, N., and Bengio, Y. (2012). Theano: new features and speed improvements. Deep Learning and Unsupervised Feature Learning NIPS 2012 Workshop. [25](#), [82](#), [212](#), [222](#), [449](#)
- Basu, S. and Christensen, J. (2013). Teaching classification boundaries to humans. In *AAAI'2013*. [329](#)
- Baxter, J. (1995). Learning internal representations. In *Proceedings of the 8th International Conference on Computational Learning Theory (COLT'95)*, pages 311–320, Santa Cruz, California. ACM Press. [246](#)
- Bayer, J. and Osendorfer, C. (2014). Learning stochastic recurrent networks. *ArXiv e-prints*. [264](#)
- Becker, S. and Hinton, G. (1992). A self-organizing neural network that discovers surfaces in random-dot stereograms. *Nature*, **355**, 161–163. [544](#)

- Behnke, S. (2001). Learning iterative image reconstruction in the neural abstraction pyramid. *Int. J. Computational Intelligence and Applications*, **1**(4), 427–438. [518](#)
- Beiu, V., Quintana, J. M., and Avedillo, M. J. (2003). VLSI implementations of threshold logic-a comprehensive survey. *Neural Networks, IEEE Transactions on*, **14**(5), 1217–1243. [454](#)
- Belkin, M. and Niyogi, P. (2002). Laplacian eigenmaps and spectral techniques for embedding and clustering. In T. Dietterich, S. Becker, and Z. Ghahramani, editors, *Advances in Neural Information Processing Systems 14 (NIPS'01)*, Cambridge, MA. MIT Press. [244](#)
- Belkin, M. and Niyogi, P. (2003). Laplacian eigenmaps for dimensionality reduction and data representation. *Neural Computation*, **15**(6), 1373–1396. [163](#), [521](#)
- Bengio, E., Bacon, P.-L., Pineau, J., and Precup, D. (2015a). Conditional computation in neural networks for faster models. arXiv:1511.06297. [453](#)
- Bengio, S. and Bengio, Y. (2000a). Taking on the curse of dimensionality in joint distributions using neural networks. *IEEE Transactions on Neural Networks, special issue on Data Mining and Knowledge Discovery*, **11**(3), 550–557. [709](#)
- Bengio, S., Vinyals, O., Jaitly, N., and Shazeer, N. (2015b). Scheduled sampling for sequence prediction with recurrent neural networks. Technical report, arXiv:1506.03099. [385](#)
- Bengio, Y. (1991). *Artificial Neural Networks and their Application to Sequence Recognition*. Ph.D. thesis, McGill University, (Computer Science), Montreal, Canada. [409](#)
- Bengio, Y. (2000). Gradient-based optimization of hyperparameters. *Neural Computation*, **12**(8), 1889–1900. [438](#)
- Bengio, Y. (2002). New distributed probabilistic language models. Technical Report 1215, Dept. IRO, Université de Montréal. [470](#)
- Bengio, Y. (2009). *Learning deep architectures for AI*. Now Publishers. [200](#), [626](#)
- Bengio, Y. (2013). Deep learning of representations: looking forward. In *Statistical Language and Speech Processing*, volume 7978 of *Lecture Notes in Computer Science*, pages 1–37. Springer, also in arXiv at <http://arxiv.org/abs/1305.0445>. [451](#)
- Bengio, Y. (2015). Early inference in energy-based models approximates back-propagation. Technical Report arXiv:1510.02777, Universite de Montreal. [658](#)
- Bengio, Y. and Bengio, S. (2000b). Modeling high-dimensional discrete data with multi-layer neural networks. In *NIPS 12*, pages 400–406. MIT Press. [707](#), [709](#), [710](#), [712](#)
- Bengio, Y. and Delalleau, O. (2009). Justifying and generalizing contrastive divergence. *Neural Computation*, **21**(6), 1601–1621. [516](#), [614](#)

- Bengio, Y. and Grandvalet, Y. (2004). No unbiased estimator of the variance of k-fold cross-validation. In S. Thrun, L. Saul, and B. Schölkopf, editors, *Advances in Neural Information Processing Systems 16 (NIPS'03)*, Cambridge, MA. MIT Press, Cambridge. [122](#)
- Bengio, Y. and LeCun, Y. (2007). Scaling learning algorithms towards AI. In *Large Scale Kernel Machines*. [19](#)
- Bengio, Y. and Monperrus, M. (2005). Non-local manifold tangent learning. In L. Saul, Y. Weiss, and L. Bottou, editors, *Advances in Neural Information Processing Systems 17 (NIPS'04)*, pages 129–136. MIT Press. [159](#), [522](#)
- Bengio, Y. and Sénecal, J.-S. (2003). Quick training of probabilistic neural nets by importance sampling. In *Proceedings of AISTATS 2003*. [473](#)
- Bengio, Y. and Sénecal, J.-S. (2008). Adaptive importance sampling to accelerate training of a neural probabilistic language model. *IEEE Trans. Neural Networks*, **19**(4), 713–722. [473](#)
- Bengio, Y., De Mori, R., Flammia, G., and Kompe, R. (1991). Phonetically motivated acoustic parameters for continuous speech recognition using artificial neural networks. In *Proceedings of EuroSpeech'91*. [27](#), [462](#)
- Bengio, Y., De Mori, R., Flammia, G., and Kompe, R. (1992). Neural network-Gaussian mixture hybrid for speech recognition or density estimation. In *NIPS 4*, pages 175–182. Morgan Kaufmann. [462](#)
- Bengio, Y., Frasconi, P., and Simard, P. (1993). The problem of learning long-term dependencies in recurrent networks. In *IEEE International Conference on Neural Networks*, pages 1183–1195, San Francisco. IEEE Press. (invited paper). [405](#)
- Bengio, Y., Simard, P., and Frasconi, P. (1994). Learning long-term dependencies with gradient descent is difficult. *IEEE Tr. Neural Nets*. [18](#), [403](#), [405](#), [406](#), [414](#)
- Bengio, Y., Latendresse, S., and Dugas, C. (1999). Gradient-based learning of hyperparameters. Learning Conference, Snowbird. [438](#)
- Bengio, Y., Ducharme, R., and Vincent, P. (2001). A neural probabilistic language model. In T. K. Leen, T. G. Dietterich, and V. Tresp, editors, *NIPS'2000*, pages 932–938. MIT Press. [18](#), [450](#), [466](#), [469](#), [475](#), [480](#), [485](#)
- Bengio, Y., Ducharme, R., Vincent, P., and Jauvin, C. (2003). A neural probabilistic language model. *JMLR*, **3**, 1137–1155. [469](#), [475](#)
- Bengio, Y., Le Roux, N., Vincent, P., Delalleau, O., and Marcotte, P. (2006a). Convex neural networks. In *NIPS'2005*, pages 123–130. [257](#)
- Bengio, Y., Delalleau, O., and Le Roux, N. (2006b). The curse of highly variable functions for local kernel machines. In *NIPS'2005*. [157](#)

- Bengio, Y., Larochelle, H., and Vincent, P. (2006c). Non-local manifold Parzen windows. In *NIPS'2005*. MIT Press. [159](#), [523](#)
- Bengio, Y., Lamblin, P., Popovici, D., and Larochelle, H. (2007). Greedy layer-wise training of deep networks. In *NIPS'2006*. [14](#), [19](#), [200](#), [324](#), [325](#), [531](#), [533](#)
- Bengio, Y., Louradour, J., Collobert, R., and Weston, J. (2009). Curriculum learning. In *ICML'09*. [329](#)
- Bengio, Y., Mesnil, G., Dauphin, Y., and Rifai, S. (2013a). Better mixing via deep representations. In *ICML'2013*. [607](#)
- Bengio, Y., Léonard, N., and Courville, A. (2013b). Estimating or propagating gradients through stochastic neurons for conditional computation. arXiv:1308.3432. [451](#), [453](#), [691](#), [693](#)
- Bengio, Y., Yao, L., Alain, G., and Vincent, P. (2013c). Generalized denoising auto-encoders as generative models. In *NIPS'2013*. [510](#), [713](#), [715](#)
- Bengio, Y., Courville, A., and Vincent, P. (2013d). Representation learning: A review and new perspectives. *IEEE Trans. Pattern Analysis and Machine Intelligence (PAMI)*, **35**(8), 1798–1828. [558](#)
- Bengio, Y., Thibodeau-Laufer, E., Alain, G., and Yosinski, J. (2014). Deep generative stochastic networks trainable by backprop. In *ICML'2014*. [713](#), [714](#), [715](#), [716](#), [717](#)
- Bennett, C. (1976). Efficient estimation of free energy differences from Monte Carlo data. *Journal of Computational Physics*, **22**(2), 245–268. [632](#)
- Bennett, J. and Lanning, S. (2007). The Netflix prize. [482](#)
- Berger, A. L., Della Pietra, V. J., and Della Pietra, S. A. (1996). A maximum entropy approach to natural language processing. *Computational Linguistics*, **22**, 39–71. [476](#)
- Berglund, M. and Raiko, T. (2013). Stochastic gradient estimate variance in contrastive divergence and persistent contrastive divergence. *CoRR*, **abs/1312.6002**. [617](#)
- Bergstra, J. (2011). *Incorporating Complex Cells into Neural Networks for Pattern Classification*. Ph.D. thesis, Université de Montréal. [254](#)
- Bergstra, J. and Bengio, Y. (2009). Slow, decorrelated features for pretraining complex cell-like networks. In *NIPS'2009*. [497](#)
- Bergstra, J. and Bengio, Y. (2012). Random search for hyper-parameter optimization. *J. Machine Learning Res.*, **13**, 281–305. [437](#), [438](#)
- Bergstra, J., Breuleux, O., Bastien, F., Lamblin, P., Pascanu, R., Desjardins, G., Turian, J., Warde-Farley, D., and Bengio, Y. (2010). Theano: a CPU and GPU math expression compiler. In *Proc. SciPy*. [25](#), [82](#), [212](#), [222](#), [449](#)

- Bergstra, J., Bardenet, R., Bengio, Y., and Kégl, B. (2011). Algorithms for hyper-parameter optimization. In *NIPS'2011*. 439
- Berkes, P. and Wiskott, L. (2005). Slow feature analysis yields a rich repertoire of complex cell properties. *Journal of Vision*, 5(6), 579–602. 498
- Bertsekas, D. P. and Tsitsiklis, J. (1996). *Neuro-Dynamic Programming*. Athena Scientific. 106
- Besag, J. (1975). Statistical analysis of non-lattice data. *The Statistician*, 24(3), 179–195. 618
- Bishop, C. M. (1994). Mixture density networks. 188
- Bishop, C. M. (1995a). Regularization and complexity control in feed-forward networks. In *Proceedings International Conference on Artificial Neural Networks ICANN'95*, volume 1, page 141–148. 242, 249
- Bishop, C. M. (1995b). Training with noise is equivalent to Tikhonov regularization. *Neural Computation*, 7(1), 108–116. 242
- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer. 98, 145
- Blum, A. L. and Rivest, R. L. (1992). Training a 3-node neural network is NP-complete. 293
- Blumer, A., Ehrenfeucht, A., Haussler, D., and Warmuth, M. K. (1989). Learnability and the Vapnik–Chervonenkis dimension. *Journal of the ACM*, 36(4), 929—865. 114
- Bonnet, G. (1964). Transformations des signaux aléatoires à travers les systèmes non linéaires sans mémoire. *Annales des Télécommunications*, 19(9–10), 203–220. 691
- Bordes, A., Weston, J., Collobert, R., and Bengio, Y. (2011). Learning structured embeddings of knowledge bases. In *AAAI 2011*. 487
- Bordes, A., Glorot, X., Weston, J., and Bengio, Y. (2012). Joint learning of words and meaning representations for open-text semantic parsing. *AISTATS'2012*. 403, 487
- Bordes, A., Glorot, X., Weston, J., and Bengio, Y. (2013a). A semantic matching energy function for learning with multi-relational data. *Machine Learning: Special Issue on Learning Semantics*. 486
- Bordes, A., Usunier, N., Garcia-Duran, A., Weston, J., and Yakhnenko, O. (2013b). Translating embeddings for modeling multi-relational data. In C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 2787–2795. Curran Associates, Inc. 487
- Bornstein, J. and Bengio, Y. (2015). Reweighted wake-sleep. In *ICLR'2015*, *arXiv:1406.2751*. 695

- Bornschein, J., Shabanian, S., Fischer, A., and Bengio, Y. (2015). Training bidirectional Helmholtz machines. Technical report, arXiv:1506.03877. [695](#)
- Boser, B. E., Guyon, I. M., and Vapnik, V. N. (1992). A training algorithm for optimal margin classifiers. In *COLT '92: Proceedings of the fifth annual workshop on Computational learning theory*, pages 144–152, New York, NY, USA. ACM. [18](#), [140](#)
- Bottou, L. (1998). Online algorithms and stochastic approximations. In D. Saad, editor, *Online Learning in Neural Networks*. Cambridge University Press, Cambridge, UK. [296](#)
- Bottou, L. (2011). From machine learning to machine reasoning. Technical report, arXiv.1102.1808. [401](#), [403](#)
- Bottou, L. (2015). Multilayer neural networks. Deep Learning Summer School. [443](#)
- Bottou, L. and Bousquet, O. (2008). The tradeoffs of large scale learning. In *NIPS'2008*. [282](#), [295](#)
- Boulanger-Lewandowski, N., Bengio, Y., and Vincent, P. (2012). Modeling temporal dependencies in high-dimensional sequences: Application to polyphonic music generation and transcription. In *ICML'12*. [688](#)
- Boureau, Y., Ponce, J., and LeCun, Y. (2010). A theoretical analysis of feature pooling in vision algorithms. In *Proc. International Conference on Machine learning (ICML'10)*. [346](#)
- Boureau, Y., Le Roux, N., Bach, F., Ponce, J., and LeCun, Y. (2011). Ask the locals: multi-way local pooling for image recognition. In *Proc. International Conference on Computer Vision (ICCV'11)*. IEEE. [346](#)
- Bourlard, H. and Kamp, Y. (1988). Auto-association by multilayer perceptrons and singular value decomposition. *Biological Cybernetics*, **59**, 291–294. [505](#)
- Bourlard, H. and Wellekens, C. (1989). Speech pattern discrimination and multi-layered perceptrons. *Computer Speech and Language*, **3**, 1–19. [462](#)
- Boyd, S. and Vandenberghe, L. (2004). *Convex Optimization*. Cambridge University Press, New York, NY, USA. [93](#)
- Brady, M. L., Raghavan, R., and Slawny, J. (1989). Back-propagation fails to separate where perceptrons succeed. *IEEE Transactions on Circuits and Systems*, **36**, 665–674. [284](#)
- Brakel, P., Stroobandt, D., and Schrauwen, B. (2013). Training energy-based models for time-series imputation. *Journal of Machine Learning Research*, **14**, 2771–2797. [676](#), [700](#)
- Brand, M. (2003). Charting a manifold. In *NIPS'2002*, pages 961–968. MIT Press. [163](#), [521](#)

- Breiman, L. (1994). Bagging predictors. *Machine Learning*, **24**(2), 123–140. [255](#)
- Breiman, L., Friedman, J. H., Olshen, R. A., and Stone, C. J. (1984). *Classification and Regression Trees*. Wadsworth International Group, Belmont, CA. [145](#)
- Bridle, J. S. (1990). Alphanets: a recurrent ‘neural’ network architecture with a hidden Markov model interpretation. *Speech Communication*, **9**(1), 83–92. [185](#)
- Briggman, K., Denk, W., Seung, S., Helmstaedter, M. N., and Turaga, S. C. (2009). Maximin affinity learning of image segmentation. In *NIPS’2009*, pages 1865–1873. [360](#)
- Brown, P. F., Cocke, J., Pietra, S. A. D., Pietra, V. J. D., Jelinek, F., Lafferty, J. D., Mercer, R. L., and Roossin, P. S. (1990). A statistical approach to machine translation. *Computational linguistics*, **16**(2), 79–85. [21](#)
- Brown, P. F., Pietra, V. J. D., DeSouza, P. V., Lai, J. C., and Mercer, R. L. (1992). Class-based  $n$ -gram models of natural language. *Computational Linguistics*, **18**, 467–479. [466](#)
- Bryson, A. and Ho, Y. (1969). *Applied optimal control: optimization, estimation, and control*. Blaisdell Pub. Co. [225](#)
- Bryson, Jr., A. E. and Denham, W. F. (1961). A steepest-ascent method for solving optimum programming problems. Technical Report BR-1303, Raytheon Company, Missle and Space Division. [225](#)
- Buciluă, C., Caruana, R., and Niculescu-Mizil, A. (2006). Model compression. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 535–541. ACM. [451](#)
- Burda, Y., Grosse, R., and Salakhutdinov, R. (2015). Importance weighted autoencoders. *arXiv preprint arXiv:1509.00519*. [700](#)
- Cai, M., Shi, Y., and Liu, J. (2013). Deep maxout neural networks for speech recognition. In *Automatic Speech Recognition and Understanding (ASRU), 2013 IEEE Workshop on*, pages 291–296. IEEE. [193](#)
- Carreira-Perpiñan, M. A. and Hinton, G. E. (2005). On contrastive divergence learning. In R. G. Cowell and Z. Ghahramani, editors, *Proceedings of the Tenth International Workshop on Artificial Intelligence and Statistics (AISTATS’05)*, pages 33–40. Society for Artificial Intelligence and Statistics. [614](#)
- Caruana, R. (1993). Multitask connectionist learning. In *Proc. 1993 Connectionist Models Summer School*, pages 372–379. [245](#)
- Cauchy, A. (1847). Méthode générale pour la résolution de systèmes d’équations simultanées. In *Compte rendu des séances de l’académie des sciences*, pages 536–538. [83](#), [224](#)

- Cayton, L. (2005). Algorithms for manifold learning. Technical Report CS2008-0923, UCSD. [163](#)
- Chandola, V., Banerjee, A., and Kumar, V. (2009). Anomaly detection: A survey. *ACM computing surveys (CSUR)*, **41**(3), 15. [102](#)
- Chapelle, O., Weston, J., and Schölkopf, B. (2003). Cluster kernels for semi-supervised learning. In S. Becker, S. Thrun, and K. Obermayer, editors, *Advances in Neural Information Processing Systems 15 (NIPS'02)*, pages 585–592, Cambridge, MA. MIT Press. [244](#)
- Chapelle, O., Schölkopf, B., and Zien, A., editors (2006). *Semi-Supervised Learning*. MIT Press, Cambridge, MA. [244](#), [544](#)
- Chellapilla, K., Puri, S., and Simard, P. (2006). High Performance Convolutional Neural Networks for Document Processing. In Guy Lorette, editor, *Tenth International Workshop on Frontiers in Handwriting Recognition*, La Baule (France). Université de Rennes 1, Suvisoft. <http://www.suvisoft.com>. [24](#), [27](#), [448](#)
- Chen, B., Ting, J.-A., Marlin, B. M., and de Freitas, N. (2010). Deep learning of invariant spatio-temporal features from video. NIPS\*2010 Deep Learning and Unsupervised Feature Learning Workshop. [361](#)
- Chen, S. F. and Goodman, J. T. (1999). An empirical study of smoothing techniques for language modeling. *Computer, Speech and Language*, **13**(4), 359–393. [465](#), [476](#)
- Chen, T., Du, Z., Sun, N., Wang, J., Wu, C., Chen, Y., and Temam, O. (2014a). DianNao: A small-footprint high-throughput accelerator for ubiquitous machine-learning. In *Proceedings of the 19th international conference on Architectural support for programming languages and operating systems*, pages 269–284. ACM. [454](#)
- Chen, T., Li, M., Li, Y., Lin, M., Wang, N., Wang, M., Xiao, T., Xu, B., Zhang, C., and Zhang, Z. (2015). MXNet: A flexible and efficient machine learning library for heterogeneous distributed systems. *arXiv preprint arXiv:1512.01274*. [25](#)
- Chen, Y., Luo, T., Liu, S., Zhang, S., He, L., Wang, J., Li, L., Chen, T., Xu, Z., Sun, N., et al. (2014b). DaDianNao: A machine-learning supercomputer. In *Microarchitecture (MICRO), 2014 47th Annual IEEE/ACM International Symposium on*, pages 609–622. IEEE. [454](#)
- Chilimbi, T., Suzue, Y., Apacible, J., and Kalyanaraman, K. (2014). Project Adam: Building an efficient and scalable deep learning training system. In *11th USENIX Symposium on Operating Systems Design and Implementation (OSDI'14)*. [450](#)
- Cho, K., Raiko, T., and Ilin, A. (2010). Parallel tempering is efficient for learning restricted Boltzmann machines. In *IJCNN'2010*. [606](#), [617](#)

- Cho, K., Raiko, T., and Ilin, A. (2011). Enhanced gradient and adaptive learning rate for training restricted Boltzmann machines. In *ICML'2011*, pages 105–112. [676](#)
- Cho, K., van Merriënboer, B., Gulcehre, C., Bougares, F., Schwenk, H., and Bengio, Y. (2014a). Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *Proceedings of the Empirical Methods in Natural Language Processing (EMNLP 2014)*. [397, 477, 478](#)
- Cho, K., Van Merriënboer, B., Bahdanau, D., and Bengio, Y. (2014b). On the properties of neural machine translation: Encoder-decoder approaches. *ArXiv e-prints*, [abs/1409.1259](#). [414](#)
- Choromanska, A., Henaff, M., Mathieu, M., Arous, G. B., and LeCun, Y. (2014). The loss surface of multilayer networks. [285, 286](#)
- Chorowski, J., Bahdanau, D., Cho, K., and Bengio, Y. (2014). End-to-end continuous speech recognition using attention-based recurrent NN: First results. arXiv:1412.1602. [463](#)
- Christianson, B. (1992). Automatic Hessians by reverse accumulation. *IMA Journal of Numerical Analysis*, **12**(2), 135–150. [224](#)
- Chrupala, G., Kadar, A., and Alishahi, A. (2015). Learning language through pictures. arXiv 1506.03694. [414](#)
- Chung, J., Gulcehre, C., Cho, K., and Bengio, Y. (2014). Empirical evaluation of gated recurrent neural networks on sequence modeling. NIPS'2014 Deep Learning workshop, arXiv 1412.3555. [414, 463](#)
- Chung, J., Gülcöhre, C., Cho, K., and Bengio, Y. (2015a). Gated feedback recurrent neural networks. In *ICML'15*. [414](#)
- Chung, J., Kastner, K., Dinh, L., Goel, K., Courville, A., and Bengio, Y. (2015b). A recurrent latent variable model for sequential data. In *NIPS'2015*. [700](#)
- Ciresan, D., Meier, U., Masci, J., and Schmidhuber, J. (2012). Multi-column deep neural network for traffic sign classification. *Neural Networks*, **32**, 333–338. [23, 200](#)
- Ciresan, D. C., Meier, U., Gambardella, L. M., and Schmidhuber, J. (2010). Deep big simple neural nets for handwritten digit recognition. *Neural Computation*, **22**, 1–14. [24, 27, 449](#)
- Coates, A. and Ng, A. Y. (2011). The importance of encoding versus training with sparse coding and vector quantization. In *ICML'2011*. [27, 254, 501](#)
- Coates, A., Lee, H., and Ng, A. Y. (2011). An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics (AISTATS 2011)*. [364, 365, 458](#)

- Coates, A., Huval, B., Wang, T., Wu, D., Catanzaro, B., and Andrew, N. (2013). Deep learning with COTS HPC systems. In S. Dasgupta and D. McAllester, editors, *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, volume 28 (3), pages 1337–1345. JMLR Workshop and Conference Proceedings. **24**, **27**, **365**, **450**
- Cohen, N., Sharir, O., and Shashua, A. (2015). On the expressive power of deep learning: A tensor analysis. arXiv:1509.05009. **557**
- Collobert, R. (2004). *Large Scale Machine Learning*. Ph.D. thesis, Université de Paris VI, LIP6. **196**
- Collobert, R. (2011). Deep learning for efficient discriminative parsing. In *AISTATS'2011*. **101**, **480**
- Collobert, R. and Weston, J. (2008a). A unified architecture for natural language processing: Deep neural networks with multitask learning. In *ICML'2008*. **474**, **480**
- Collobert, R. and Weston, J. (2008b). A unified architecture for natural language processing: Deep neural networks with multitask learning. In *ICML'2008*. **538**
- Collobert, R., Bengio, S., and Bengio, Y. (2001). A parallel mixture of SVMs for very large scale problems. Technical Report IDIAP-RR-01-12, IDIAP. **453**
- Collobert, R., Bengio, S., and Bengio, Y. (2002). Parallel mixture of SVMs for very large scale problems. *Neural Computation*, **14**(5), 1105–1114. **453**
- Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., and Kuksa, P. (2011a). Natural language processing (almost) from scratch. *The Journal of Machine Learning Research*, **12**, 2493–2537. **329**, **480**, **538**, **539**
- Collobert, R., Kavukcuoglu, K., and Farabet, C. (2011b). Torch7: A Matlab-like environment for machine learning. In *BigLearn, NIPS Workshop*. **25**, **210**, **449**
- Comon, P. (1994). Independent component analysis - a new concept? *Signal Processing*, **36**, 287–314. **494**
- Cortes, C. and Vapnik, V. (1995). Support vector networks. *Machine Learning*, **20**, 273–297. **18**, **140**
- Couprise, C., Farabet, C., Najman, L., and LeCun, Y. (2013). Indoor semantic segmentation using depth information. In *International Conference on Learning Representations (ICLR2013)*. **23**, **200**
- Courbariaux, M., Bengio, Y., and David, J.-P. (2015). Low precision arithmetic for deep learning. In *Arxiv:1412.7024, ICLR'2015 Workshop*. **455**
- Courville, A., Bergstra, J., and Bengio, Y. (2011). Unsupervised models of images by spike-and-slab RBMs. In *ICML'11*. **564**, **683**

- Courville, A., Desjardins, G., Bergstra, J., and Bengio, Y. (2014). The spike-and-slab RBM and extensions to discrete and sparse data distributions. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, **36**(9), 1874–1887. [685](#)
- Cover, T. M. and Thomas, J. A. (2006). *Elements of Information Theory, 2nd Edition*. Wiley-Interscience. [73](#)
- Cox, D. and Pinto, N. (2011). Beyond simple features: A large-scale feature search approach to unconstrained face recognition. In *Automatic Face & Gesture Recognition and Workshops (FG 2011), 2011 IEEE International Conference on*, pages 8–15. IEEE. [364](#)
- Cramér, H. (1946). *Mathematical methods of statistics*. Princeton University Press. [135](#), [295](#)
- Crick, F. H. C. and Mitchison, G. (1983). The function of dream sleep. *Nature*, **304**, 111–114. [612](#)
- Cybenko, G. (1989). Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals, and Systems*, **2**, 303–314. [197](#)
- Dahl, G. E., Ranzato, M., Mohamed, A., and Hinton, G. E. (2010). Phone recognition with the mean-covariance restricted Boltzmann machine. In *NIPS’2010*. [23](#)
- Dahl, G. E., Yu, D., Deng, L., and Acero, A. (2012). Context-dependent pre-trained deep neural networks for large vocabulary speech recognition. *IEEE Transactions on Audio, Speech, and Language Processing*, **20**(1), 33–42. [462](#)
- Dahl, G. E., Sainath, T. N., and Hinton, G. E. (2013). Improving deep neural networks for LVCSR using rectified linear units and dropout. In *ICASSP’2013*. [462](#)
- Dahl, G. E., Jaitly, N., and Salakhutdinov, R. (2014). Multi-task neural networks for QSAR predictions. arXiv:1406.1231. [26](#)
- Dauphin, Y. and Bengio, Y. (2013). Stochastic ratio matching of RBMs for sparse high-dimensional inputs. In *NIPS26*. NIPS Foundation. [622](#)
- Dauphin, Y., Glorot, X., and Bengio, Y. (2011). Large-scale learning of embeddings with reconstruction sampling. In *ICML’2011*. [474](#)
- Dauphin, Y., Pascanu, R., Gulcehre, C., Cho, K., Ganguli, S., and Bengio, Y. (2014). Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. In *NIPS’2014*. [285](#), [286](#), [288](#)
- Davis, A., Rubinstein, M., Wadhwa, N., Mysore, G., Durand, F., and Freeman, W. T. (2014). The visual microphone: Passive recovery of sound from video. *ACM Transactions on Graphics (Proc. SIGGRAPH)*, **33**(4), 79:1–79:10. [455](#)

- Dayan, P. (1990). Reinforcement comparison. In *Connectionist Models: Proceedings of the 1990 Connectionist Summer School*, San Mateo, CA. 693
- Dayan, P. and Hinton, G. E. (1996). Varieties of Helmholtz machine. *Neural Networks*, **9**(8), 1385–1403. 695
- Dayan, P., Hinton, G. E., Neal, R. M., and Zemel, R. S. (1995). The Helmholtz machine. *Neural computation*, **7**(5), 889–904. 695
- Dean, J., Corrado, G., Monga, R., Chen, K., Devin, M., Le, Q., Mao, M., Ranzato, M., Senior, A., Tucker, P., Yang, K., and Ng, A. Y. (2012). Large scale distributed deep networks. In *NIPS'2012*. 25, 450
- Dean, T. and Kanazawa, K. (1989). A model for reasoning about persistence and causation. *Computational Intelligence*, **5**(3), 142–150. 664
- Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K., and Harshman, R. (1990). Indexing by latent semantic analysis. *Journal of the American Society for Information Science*, **41**(6), 391–407. 479, 485
- Delalleau, O. and Bengio, Y. (2011). Shallow vs. deep sum-product networks. In *NIPS*. 19, 557
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. (2009). ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*. 21
- Deng, J., Berg, A. C., Li, K., and Fei-Fei, L. (2010a). What does classifying more than 10,000 image categories tell us? In *Proceedings of the 11th European Conference on Computer Vision: Part V*, ECCV'10, pages 71–84, Berlin, Heidelberg. Springer-Verlag. 21
- Deng, L. and Yu, D. (2014). Deep learning – methods and applications. *Foundations and Trends in Signal Processing*. 463
- Deng, L., Seltzer, M., Yu, D., Acero, A., Mohamed, A., and Hinton, G. (2010b). Binary coding of speech spectrograms using a deep auto-encoder. In *Interspeech 2010*, Makuhari, Chiba, Japan. 23
- Denil, M., Bazzani, L., Larochelle, H., and de Freitas, N. (2012). Learning where to attend with deep architectures for image tracking. *Neural Computation*, **24**(8), 2151–2184. 368
- Denton, E., Chintala, S., Szlam, A., and Fergus, R. (2015). Deep generative image models using a Laplacian pyramid of adversarial networks. *NIPS*. 703, 704, 720
- Desjardins, G. and Bengio, Y. (2008). Empirical evaluation of convolutional RBMs for vision. Technical Report 1327, Département d’Informatique et de Recherche Opérationnelle, Université de Montréal. 685

- Desjardins, G., Courville, A. C., Bengio, Y., Vincent, P., and Delalleau, O. (2010). Tempered Markov chain Monte Carlo for training of restricted Boltzmann machines. In *International Conference on Artificial Intelligence and Statistics*, pages 145–152. **606**, **617**
- Desjardins, G., Courville, A., and Bengio, Y. (2011). On tracking the partition function. In *NIPS'2011*. **633**
- Desjardins, G., Simonyan, K., Pascanu, R., *et al.* (2015). Natural neural networks. In *Advances in Neural Information Processing Systems*, pages 2062–2070. **321**
- Devlin, J., Zbib, R., Huang, Z., Lamar, T., Schwartz, R., and Makhoul, J. (2014). Fast and robust neural network joint models for statistical machine translation. In *Proc. ACL'2014*. **476**
- Devroye, L. (2013). *Non-Uniform Random Variate Generation*. SpringerLink : Bücher. Springer New York. **696**
- DiCarlo, J. J. (2013). Mechanisms underlying visual object recognition: Humans vs. neurons vs. machines. *NIPS Tutorial*. **26**, **367**
- Dinh, L., Krueger, D., and Bengio, Y. (2014). NICE: Non-linear independent components estimation. arXiv:1410.8516. **496**
- Donahue, J., Hendricks, L. A., Guadarrama, S., Rohrbach, M., Venugopalan, S., Saenko, K., and Darrell, T. (2014). Long-term recurrent convolutional networks for visual recognition and description. arXiv:1411.4389. **102**
- Donoho, D. L. and Grimes, C. (2003). Hessian eigenmaps: new locally linear embedding techniques for high-dimensional data. Technical Report 2003-08, Dept. Statistics, Stanford University. **163**, **522**
- Dosovitskiy, A., Springenberg, J. T., and Brox, T. (2015). Learning to generate chairs with convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1538–1546. **697**, **706**, **707**
- Doya, K. (1993). Bifurcations of recurrent neural networks in gradient descent learning. *IEEE Transactions on Neural Networks*, **1**, 75–80. **403**, **406**
- Dreyfus, S. E. (1962). The numerical solution of variational problems. *Journal of Mathematical Analysis and Applications*, **5**(**1**), 30–45. **225**
- Dreyfus, S. E. (1973). The computational solution of optimal control problems with time lag. *IEEE Transactions on Automatic Control*, **18**(**4**), 383–385. **225**
- Drucker, H. and LeCun, Y. (1992). Improving generalisation performance using double back-propagation. *IEEE Transactions on Neural Networks*, **3**(**6**), 991–997. **270**

- Duchi, J., Hazan, E., and Singer, Y. (2011). Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*. **307**
- Dudik, M., Langford, J., and Li, L. (2011). Doubly robust policy evaluation and learning. In *Proceedings of the 28th International Conference on Machine learning*, ICML '11. **485**
- Dugas, C., Bengio, Y., Bélisle, F., and Nadeau, C. (2001). Incorporating second-order functional knowledge for better option pricing. In T. Leen, T. Dietterich, and V. Tresp, editors, *Advances in Neural Information Processing Systems 13 (NIPS'00)*, pages 472–478. MIT Press. **68, 196**
- Dziugaite, G. K., Roy, D. M., and Ghahramani, Z. (2015). Training generative neural networks via maximum mean discrepancy optimization. *arXiv preprint arXiv:1505.03906*. **705**
- El Hihi, S. and Bengio, Y. (1996). Hierarchical recurrent neural networks for long-term dependencies. In *NIPS'1995*. **401, 410**
- Elkahky, A. M., Song, Y., and He, X. (2015). A multi-view deep learning approach for cross domain user modeling in recommendation systems. In *Proceedings of the 24th International Conference on World Wide Web*, pages 278–288. **488**
- Elman, J. L. (1993). Learning and development in neural networks: The importance of starting small. *Cognition*, **48**, 781–799. **329**
- Erhan, D., Manzagol, P.-A., Bengio, Y., Bengio, S., and Vincent, P. (2009). The difficulty of training deep architectures and the effect of unsupervised pre-training. In *Proceedings of AISTATS'2009*. **200**
- Erhan, D., Bengio, Y., Courville, A., Manzagol, P., Vincent, P., and Bengio, S. (2010). Why does unsupervised pre-training help deep learning? *J. Machine Learning Res.* **532, 536, 537**
- Fahlman, S. E., Hinton, G. E., and Sejnowski, T. J. (1983). Massively parallel architectures for AI: NETL, thistle, and Boltzmann machines. In *Proceedings of the National Conference on Artificial Intelligence AAAI-83*. **573, 656**
- Fang, H., Gupta, S., Iandola, F., Srivastava, R., Deng, L., Dollár, P., Gao, J., He, X., Mitchell, M., Platt, J. C., Zitnick, C. L., and Zweig, G. (2015). From captions to visual concepts and back. *arXiv:1411.4952*. **102**
- Farabet, C., LeCun, Y., Kavukcuoglu, K., Culurciello, E., Martini, B., Akselrod, P., and Talay, S. (2011). Large-scale FPGA-based convolutional networks. In R. Bekkerman, M. Bilenko, and J. Langford, editors, *Scaling up Machine Learning: Parallel and Distributed Approaches*. Cambridge University Press. **526**

- Farabet, C., Couprie, C., Najman, L., and LeCun, Y. (2013). Learning hierarchical features for scene labeling. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **35**(8), 1915–1929. [23](#), [200](#), [360](#)
- Fei-Fei, L., Fergus, R., and Perona, P. (2006). One-shot learning of object categories. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **28**(4), 594–611. [541](#)
- Finn, C., Tan, X. Y., Duan, Y., Darrell, T., Levine, S., and Abbeel, P. (2015). Learning visual feature spaces for robotic manipulation with deep spatial autoencoders. *arXiv preprint arXiv:1509.06113*. [25](#)
- Fisher, R. A. (1936). The use of multiple measurements in taxonomic problems. *Annals of Eugenics*, **7**, 179–188. [21](#), [105](#)
- Földiák, P. (1989). Adaptive network for optimal linear feature extraction. In *International Joint Conference on Neural Networks (IJCNN)*, volume 1, pages 401–405, Washington 1989. IEEE, New York. [497](#)
- Franzius, M., Sprekeler, H., and Wiskott, L. (2007). Slowness and sparseness lead to place, head-direction, and spatial-view cells. [498](#)
- Franzius, M., Wilbert, N., and Wiskott, L. (2008). Invariant object recognition with slow feature analysis. In *Artificial Neural Networks-ICANN 2008*, pages 961–970. Springer. [499](#)
- Frasconi, P., Gori, M., and Sperduti, A. (1997). On the efficient classification of data structures by neural networks. In *Proc. Int. Joint Conf. on Artificial Intelligence*. [401](#), [403](#)
- Frasconi, P., Gori, M., and Sperduti, A. (1998). A general framework for adaptive processing of data structures. *IEEE Transactions on Neural Networks*, **9**(5), 768–786. [401](#), [403](#)
- Freund, Y. and Schapire, R. E. (1996a). Experiments with a new boosting algorithm. In *Machine Learning: Proceedings of Thirteenth International Conference*, pages 148–156, USA. ACM. [257](#)
- Freund, Y. and Schapire, R. E. (1996b). Game theory, on-line prediction and boosting. In *Proceedings of the Ninth Annual Conference on Computational Learning Theory*, pages 325–332. [257](#)
- Frey, B. J. (1998). *Graphical models for machine learning and digital communication*. MIT Press. [707](#), [708](#)
- Frey, B. J., Hinton, G. E., and Dayan, P. (1996). Does the wake-sleep algorithm learn good density estimators? In D. Touretzky, M. Mozer, and M. Hasselmo, editors, *Advances in Neural Information Processing Systems 8 (NIPS'95)*, pages 661–670. MIT Press, Cambridge, MA. [654](#)

- Frobenius, G. (1908). Über matrizen aus positiven elementen, s. *B. Preuss. Akad. Wiss. Berlin, Germany.* 600
- Fukushima, K. (1975). Cognitron: A self-organizing multilayered neural network. *Biological Cybernetics*, **20**, 121–136. 16, 226, 531
- Fukushima, K. (1980). Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position. *Biological Cybernetics*, **36**, 193–202. 16, 24, 27, 226, 368
- Gal, Y. and Ghahramani, Z. (2015). Bayesian convolutional neural networks with Bernoulli approximate variational inference. *arXiv preprint arXiv:1506.02158*. 263
- Gallinari, P., LeCun, Y., Thiria, S., and Fogelman-Soulie, F. (1987). Memoires associatives distribuees. In *Proceedings of COGNITIVA 87*, Paris, La Villette. 518
- Garcia-Duran, A., Bordes, A., Usunier, N., and Grandvalet, Y. (2015). Combining two and three-way embeddings models for link prediction in knowledge bases. *arXiv preprint arXiv:1506.00999*. 487
- Garofolo, J. S., Lamel, L. F., Fisher, W. M., Fiscus, J. G., and Pallett, D. S. (1993). Darpa timit acoustic-phonetic continuous speech corpus cd-rom. nist speech disc 1-1.1. *NASA STI/Recon Technical Report N*, **93**, 27403. 462
- Garson, J. (1900). The metric system of identification of criminals, as used in Great Britain and Ireland. *The Journal of the Anthropological Institute of Great Britain and Ireland*, (2), 177–227. 21
- Gers, F. A., Schmidhuber, J., and Cummins, F. (2000). Learning to forget: Continual prediction with LSTM. *Neural computation*, **12**(10), 2451–2471. 411, 415
- Ghahramani, Z. and Hinton, G. E. (1996). The EM algorithm for mixtures of factor analyzers. Technical Report CRG-TR-96-1, Dpt. of Comp. Sci., Univ. of Toronto. 492
- Gillick, D., Brunk, C., Vinyals, O., and Subramanya, A. (2015). Multilingual language processing from bytes. *arXiv preprint arXiv:1512.00103*. 480
- Girshick, R., Donahue, J., Darrell, T., and Malik, J. (2015). Region-based convolutional networks for accurate object detection and segmentation. 429
- Giudice, M. D., Manera, V., and Keysers, C. (2009). Programmed to learn? The ontogeny of mirror neurons. *Dev. Sci.*, **12**(2), 350—363. 658
- Glorot, X. and Bengio, Y. (2010). Understanding the difficulty of training deep feedforward neural networks. In *AISTATS'2010*. 303
- Glorot, X., Bordes, A., and Bengio, Y. (2011a). Deep sparse rectifier neural networks. In *AISTATS'2011*. 16, 173, 196, 226

- Glorot, X., Bordes, A., and Bengio, Y. (2011b). Domain adaptation for large-scale sentiment classification: A deep learning approach. In *ICML'2011*. 510, 540
- Goldberger, J., Roweis, S., Hinton, G. E., and Salakhutdinov, R. (2005). Neighbourhood components analysis. In L. Saul, Y. Weiss, and L. Bottou, editors, *Advances in Neural Information Processing Systems 17 (NIPS'04)*. MIT Press. 115
- Gong, S., McKenna, S., and Psarrou, A. (2000). *Dynamic Vision: From Images to Face Recognition*. Imperial College Press. 164, 522
- Goodfellow, I., Le, Q., Saxe, A., and Ng, A. (2009). Measuring invariances in deep networks. In *NIPS'2009*, pages 646–654. 254
- Goodfellow, I., Koenig, N., Muja, M., Pantofaru, C., Sorokin, A., and Takayama, L. (2010). Help me help you: Interfaces for personal robots. In *Proc. of Human Robot Interaction (HRI)*, Osaka, Japan. ACM Press, ACM Press. 100
- Goodfellow, I. J. (2010). Technical report: Multidimensional, downsampled convolution for autoencoders. Technical report, Université de Montréal. 358
- Goodfellow, I. J. (2014). On distinguishability criteria for estimating generative models. In *International Conference on Learning Representations, Workshops Track*. 625, 702, 703
- Goodfellow, I. J., Courville, A., and Bengio, Y. (2011). Spike-and-slab sparse coding for unsupervised feature discovery. In *NIPS Workshop on Challenges in Learning Hierarchical Models*. 535, 541
- Goodfellow, I. J., Warde-Farley, D., Mirza, M., Courville, A., and Bengio, Y. (2013a). Maxout networks. In S. Dasgupta and D. McAllester, editors, *ICML'13*, pages 1319–1327. 192, 263, 345, 366, 458
- Goodfellow, I. J., Mirza, M., Courville, A., and Bengio, Y. (2013b). Multi-prediction deep Boltzmann machines. In *NIPS'26*. NIPS Foundation. 100, 620, 673, 674, 675, 676, 677, 700
- Goodfellow, I. J., Warde-Farley, D., Lamblin, P., Dumoulin, V., Mirza, M., Pascanu, R., Bergstra, J., Bastien, F., and Bengio, Y. (2013c). Pylearn2: a machine learning research library. *arXiv preprint arXiv:1308.4214*. 25, 449
- Goodfellow, I. J., Courville, A., and Bengio, Y. (2013d). Scaling up spike-and-slab models for unsupervised feature learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **35**(8), 1902–1914. 500, 501, 502, 652, 685
- Goodfellow, I. J., Mirza, M., Xiao, D., Courville, A., and Bengio, Y. (2014a). An empirical investigation of catastrophic forgetting in gradient-based neural networks. In *ICLR'2014*. 193

- Goodfellow, I. J., Shlens, J., and Szegedy, C. (2014b). Explaining and harnessing adversarial examples. *CoRR*, **abs/1412.6572**. 267, 268, 270, 558, 559
- Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014c). Generative adversarial networks. In *NIPS'2014*. 547, 691, 702, 703, 706
- Goodfellow, I. J., Bulatov, Y., Ibarz, J., Arnoud, S., and Shet, V. (2014d). Multi-digit number recognition from Street View imagery using deep convolutional neural networks. In *International Conference on Learning Representations*. 25, 101, 200, 201, 202, 391, 425, 452
- Goodfellow, I. J., Vinyals, O., and Saxe, A. M. (2015). Qualitatively characterizing neural network optimization problems. In *International Conference on Learning Representations*. 285, 286, 287, 291
- Goodman, J. (2001). Classes for fast maximum entropy training. In *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Utah. 470
- Gori, M. and Tesi, A. (1992). On the problem of local minima in backpropagation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **PAMI-14**(1), 76–86. 284
- Gosset, W. S. (1908). The probable error of a mean. *Biometrika*, **6**(1), 1–25. Originally published under the pseudonym “Student”. 21
- Gouws, S., Bengio, Y., and Corrado, G. (2014). BilBOWA: Fast bilingual distributed representations without word alignments. Technical report, arXiv:1410.2455. 479, 542
- Graf, H. P. and Jackel, L. D. (1989). Analog electronic neural network circuits. *Circuits and Devices Magazine, IEEE*, **5**(4), 44–49. 454
- Graves, A. (2011). Practical variational inference for neural networks. In *NIPS'2011*. 242
- Graves, A. (2012). *Supervised Sequence Labelling with Recurrent Neural Networks*. Studies in Computational Intelligence. Springer. 375, 396, 414, 463
- Graves, A. (2013). Generating sequences with recurrent neural networks. Technical report, arXiv:1308.0850. 189, 411, 418, 422
- Graves, A. and Jaitly, N. (2014). Towards end-to-end speech recognition with recurrent neural networks. In *ICML'2014*. 411
- Graves, A. and Schmidhuber, J. (2005). Framewise phoneme classification with bidirectional LSTM and other neural network architectures. *Neural Networks*, **18**(5), 602–610. 396
- Graves, A. and Schmidhuber, J. (2009). Offline handwriting recognition with multidimensional recurrent neural networks. In D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, editors, *NIPS'2008*, pages 545–552. 396

- Graves, A., Fernández, S., Gomez, F., and Schmidhuber, J. (2006). Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks. In *ICML'2006*, pages 369–376, Pittsburgh, USA. **463**
- Graves, A., Liwicki, M., Bunke, H., Schmidhuber, J., and Fernández, S. (2008). Unconstrained on-line handwriting recognition with recurrent neural networks. In J. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *NIPS'2007*, pages 577–584. **396**
- Graves, A., Liwicki, M., Fernández, S., Bertolami, R., Bunke, H., and Schmidhuber, J. (2009). A novel connectionist system for unconstrained handwriting recognition. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, **31**(5), 855–868. **411**
- Graves, A., Mohamed, A., and Hinton, G. (2013). Speech recognition with deep recurrent neural networks. In *ICASSP'2013*, pages 6645–6649. **396, 399, 401, 411, 413, 414, 463**
- Graves, A., Wayne, G., and Danihelka, I. (2014a). Neural Turing machines. arXiv:1410.5401. **25**
- Graves, A., Wayne, G., and Danihelka, I. (2014b). Neural Turing machines. *arXiv preprint arXiv:1410.5401*. **419, 421**
- Grefenstette, E., Hermann, K. M., Suleyman, M., and Blunsom, P. (2015). Learning to transduce with unbounded memory. In *NIPS'2015*. **421**
- Greff, K., Srivastava, R. K., Koutník, J., Steunebrink, B. R., and Schmidhuber, J. (2015). LSTM: a search space odyssey. *arXiv preprint arXiv:1503.04069*. **415**
- Gregor, K. and LeCun, Y. (2010a). Emergence of complex-like cells in a temporal product network with local receptive fields. Technical report, arXiv:1006.0448. **353**
- Gregor, K. and LeCun, Y. (2010b). Learning fast approximations of sparse coding. In L. Bottou and M. Littman, editors, *Proceedings of the Twenty-seventh International Conference on Machine Learning (ICML-10)*. ACM. **655**
- Gregor, K., Danihelka, I., Mnih, A., Blundell, C., and Wierstra, D. (2014). Deep autoregressive networks. In *International Conference on Machine Learning (ICML'2014)*. **695**
- Gregor, K., Danihelka, I., Graves, A., and Wierstra, D. (2015). DRAW: A recurrent neural network for image generation. *arXiv preprint arXiv:1502.04623*. **700**
- Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., and Smola, A. (2012). A kernel two-sample test. *The Journal of Machine Learning Research*, **13**(1), 723–773. **705**
- Gülçehre, Ç. and Bengio, Y. (2013). Knowledge matters: Importance of prior information for optimization. In *International Conference on Learning Representations (ICLR'2013)*. **25**

- Guo, H. and Gelfand, S. B. (1992). Classification trees with neural network feature extraction. *Neural Networks, IEEE Transactions on*, **3**(6), 923–933. [453](#)
- Gupta, S., Agrawal, A., Gopalakrishnan, K., and Narayanan, P. (2015). Deep learning with limited numerical precision. *CoRR*, **abs/1502.02551**. [455](#)
- Gutmann, M. and Hyvarinen, A. (2010). Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of The Thirteenth International Conference on Artificial Intelligence and Statistics (AISTATS'10)*. [623](#)
- Hadsell, R., Sermanet, P., Ben, J., Erkan, A., Han, J., Muller, U., and LeCun, Y. (2007). Online learning for offroad robots: Spatial label propagation to learn long-range traversability. In *Proceedings of Robotics: Science and Systems*, Atlanta, GA, USA. [456](#)
- Hajnal, A., Maass, W., Pudlak, P., Szegedy, M., and Turan, G. (1993). Threshold circuits of bounded depth. *J. Comput. System. Sci.*, **46**, 129–154. [198](#)
- Håstad, J. (1986). Almost optimal lower bounds for small depth circuits. In *Proceedings of the 18th annual ACM Symposium on Theory of Computing*, pages 6–20, Berkeley, California. ACM Press. [198](#)
- Håstad, J. and Goldmann, M. (1991). On the power of small-depth threshold circuits. *Computational Complexity*, **1**, 113–129. [198](#)
- Hastie, T., Tibshirani, R., and Friedman, J. (2001). *The elements of statistical learning: data mining, inference and prediction*. Springer Series in Statistics. Springer Verlag. [145](#)
- He, K., Zhang, X., Ren, S., and Sun, J. (2015). Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification. *arXiv preprint arXiv:1502.01852*. [28, 192](#)
- Hebb, D. O. (1949). *The Organization of Behavior*. Wiley, New York. [14, 17, 658](#)
- Henaff, M., Jarrett, K., Kavukcuoglu, K., and LeCun, Y. (2011). Unsupervised learning of sparse features for scalable audio classification. In *ISMIR'11*. [526](#)
- Henderson, J. (2003). Inducing history representations for broad coverage statistical parsing. In *HLT-NAACL*, pages 103–110. [480](#)
- Henderson, J. (2004). Discriminative training of a neural network statistical parser. In *Proceedings of the 42nd Annual Meeting on Association for Computational Linguistics*, page 95. [480](#)
- Henniges, M., Puertas, G., Bornschein, J., Eggert, J., and Lücke, J. (2010). Binary sparse coding. In *Latent Variable Analysis and Signal Separation*, pages 450–457. Springer. [643](#)

- Herault, J. and Ans, B. (1984). Circuits neuronaux à synapses modifiables: Décodage de messages composites par apprentissage non supervisé. *Comptes Rendus de l'Académie des Sciences*, **299(III-13)**, 525—528. [494](#)
- Hinton, G. (2012). Neural networks for machine learning. Coursera, video lectures. [307](#)
- Hinton, G., Deng, L., Dahl, G. E., Mohamed, A., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Sainath, T., and Kingsbury, B. (2012a). Deep neural networks for acoustic modeling in speech recognition. *IEEE Signal Processing Magazine*, **29**(6), 82–97. [23](#), [463](#)
- Hinton, G., Vinyals, O., and Dean, J. (2015). Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*. [451](#)
- Hinton, G. E. (1989). Connectionist learning procedures. *Artificial Intelligence*, **40**, 185–234. [497](#)
- Hinton, G. E. (1990). Mapping part-whole hierarchies into connectionist networks. *Artificial Intelligence*, **46**(1), 47–75. [421](#)
- Hinton, G. E. (1999). Products of experts. In *ICANN'1999*. [573](#)
- Hinton, G. E. (2000). Training products of experts by minimizing contrastive divergence. Technical Report GCNU TR 2000-004, Gatsby Unit, University College London. [613](#), [678](#)
- Hinton, G. E. (2006). To recognize shapes, first learn to generate images. Technical Report UTML TR 2006-003, University of Toronto. [531](#), [598](#)
- Hinton, G. E. (2007a). How to do backpropagation in a brain. Invited talk at the NIPS'2007 Deep Learning Workshop. [658](#)
- Hinton, G. E. (2007b). Learning multiple layers of representation. *Trends in cognitive sciences*, **11**(10), 428–434. [662](#)
- Hinton, G. E. (2010). A practical guide to training restricted Boltzmann machines. Technical Report UTML TR 2010-003, Department of Computer Science, University of Toronto. [613](#)
- Hinton, G. E. and Ghahramani, Z. (1997). Generative models for discovering sparse distributed representations. *Philosophical Transactions of the Royal Society of London*. [146](#)
- Hinton, G. E. and McClelland, J. L. (1988). Learning representations by recirculation. In *NIPS'1987*, pages 358–366. [505](#)
- Hinton, G. E. and Roweis, S. (2003). Stochastic neighbor embedding. In *NIPS'2002*. [522](#)

- Hinton, G. E. and Salakhutdinov, R. (2006). Reducing the dimensionality of data with neural networks. *Science*, **313**(5786), 504–507. [512](#), [527](#), [531](#), [532](#), [537](#)
- Hinton, G. E. and Sejnowski, T. J. (1986). Learning and relearning in Boltzmann machines. In D. E. Rumelhart and J. L. McClelland, editors, *Parallel Distributed Processing*, volume 1, chapter 7, pages 282–317. MIT Press, Cambridge. [573](#), [656](#)
- Hinton, G. E. and Sejnowski, T. J. (1999). *Unsupervised learning: foundations of neural computation*. MIT press. [544](#)
- Hinton, G. E. and Shallice, T. (1991). Lesioning an attractor network: investigations of acquired dyslexia. *Psychological review*, **98**(1), 74. [13](#)
- Hinton, G. E. and Zemel, R. S. (1994). Autoencoders, minimum description length, and Helmholtz free energy. In *NIPS'93*. [505](#)
- Hinton, G. E., Sejnowski, T. J., and Ackley, D. H. (1984). Boltzmann machines: Constraint satisfaction networks that learn. Technical Report TR-CMU-CS-84-119, Carnegie-Mellon University, Dept. of Computer Science. [573](#), [656](#)
- Hinton, G. E., McClelland, J., and Rumelhart, D. (1986). Distributed representations. In D. E. Rumelhart and J. L. McClelland, editors, *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, volume 1, pages 77–109. MIT Press, Cambridge. [17](#), [225](#), [529](#)
- Hinton, G. E., Revow, M., and Dayan, P. (1995a). Recognizing handwritten digits using mixtures of linear models. In G. Tesauro, D. Touretzky, and T. Leen, editors, *Advances in Neural Information Processing Systems 7 (NIPS'94)*, pages 1015–1022. MIT Press, Cambridge, MA. [492](#)
- Hinton, G. E., Dayan, P., Frey, B. J., and Neal, R. M. (1995b). The wake-sleep algorithm for unsupervised neural networks. *Science*, **268**, 1558–1161. [507](#), [654](#)
- Hinton, G. E., Dayan, P., and Revow, M. (1997). Modelling the manifolds of images of handwritten digits. *IEEE Transactions on Neural Networks*, **8**, 65–74. [502](#)
- Hinton, G. E., Welling, M., Teh, Y. W., and Osindero, S. (2001). A new view of ICA. In *Proceedings of 3rd International Conference on Independent Component Analysis and Blind Signal Separation (ICA '01)*, pages 746–751, San Diego, CA. [494](#)
- Hinton, G. E., Osindero, S., and Teh, Y. (2006). A fast learning algorithm for deep belief nets. *Neural Computation*, **18**, 1527–1554. [14](#), [19](#), [27](#), [142](#), [531](#), [532](#), [662](#), [663](#)
- Hinton, G. E., Deng, L., Yu, D., Dahl, G. E., Mohamed, A., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Sainath, T. N., and Kingsbury, B. (2012b). Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal Process. Mag.*, **29**(6), 82–97. [101](#)

- Hinton, G. E., Srivastava, N., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2012c). Improving neural networks by preventing co-adaptation of feature detectors. Technical report, arXiv:1207.0580. 239, 261, 266
- Hinton, G. E., Vinyals, O., and Dean, J. (2014). Dark knowledge. Invited talk at the BayLearn Bay Area Machine Learning Symposium. 451
- Hochreiter, S. (1991). Untersuchungen zu dynamischen neuronalen Netzen. Diploma thesis, T.U. München. 18, 403, 405
- Hochreiter, S. and Schmidhuber, J. (1995). Simplifying neural nets by discovering flat minima. In *Advances in Neural Information Processing Systems 7*, pages 529–536. MIT Press. 243
- Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735–1780. 18, 411, 414
- Hochreiter, S., Informatik, F. F., Bengio, Y., Frasconi, P., and Schmidhuber, J. (2000). Gradient flow in recurrent nets: the difficulty of learning long-term dependencies. In J. Kolen and S. Kremer, editors, *Field Guide to Dynamical Recurrent Networks*. IEEE Press. 414
- Holi, J. L. and Hwang, J.-N. (1993). Finite precision error analysis of neural network hardware implementations. *Computers, IEEE Transactions on*, 42(3), 281–290. 454
- Holt, J. L. and Baker, T. E. (1991). Back propagation simulations using limited precision calculations. In *Neural Networks, 1991., IJCNN-91-Seattle International Joint Conference on*, volume 2, pages 121–126. IEEE. 454
- Hornik, K., Stinchcombe, M., and White, H. (1989). Multilayer feedforward networks are universal approximators. *Neural Networks*, 2, 359–366. 197
- Hornik, K., Stinchcombe, M., and White, H. (1990). Universal approximation of an unknown mapping and its derivatives using multilayer feedforward networks. *Neural networks*, 3(5), 551–560. 197
- Hsu, F.-H. (2002). *Behind Deep Blue: Building the Computer That Defeated the World Chess Champion*. Princeton University Press, Princeton, NJ, USA. 2
- Huang, F. and Ogata, Y. (2002). Generalized pseudo-likelihood estimates for Markov random fields on lattice. *Annals of the Institute of Statistical Mathematics*, 54(1), 1–18. 619
- Huang, P.-S., He, X., Gao, J., Deng, L., Acero, A., and Heck, L. (2013). Learning deep structured semantic models for web search using clickthrough data. In *Proceedings of the 22nd ACM international conference on Conference on information & knowledge management*, pages 2333–2338. ACM. 483

- Hubel, D. and Wiesel, T. (1968). Receptive fields and functional architecture of monkey striate cortex. *Journal of Physiology (London)*, **195**, 215–243. [365](#)
- Hubel, D. H. and Wiesel, T. N. (1959). Receptive fields of single neurons in the cat's striate cortex. *Journal of Physiology*, **148**, 574–591. [365](#)
- Hubel, D. H. and Wiesel, T. N. (1962). Receptive fields, binocular interaction, and functional architecture in the cat's visual cortex. *Journal of Physiology (London)*, **160**, 106–154. [365](#)
- Huszar, F. (2015). How (not) to train your generative model: schedule sampling, likelihood, adversary? *arXiv:1511.05101*. [699](#)
- Hutter, F., Hoos, H., and Leyton-Brown, K. (2011). Sequential model-based optimization for general algorithm configuration. In *LION-5*. Extended version as UBC Tech report TR-2010-10. [439](#)
- Hyotyniemi, H. (1996). Turing machines are recurrent neural networks. In *STeP'96*, pages 13–24. [380](#)
- Hyvärinen, A. (1999). Survey on independent component analysis. *Neural Computing Surveys*, **2**, 94–128. [494](#)
- Hyvärinen, A. (2005). Estimation of non-normalized statistical models using score matching. *Journal of Machine Learning Research*, **6**, 695–709. [516](#), [620](#)
- Hyvärinen, A. (2007a). Connections between score matching, contrastive divergence, and pseudolikelihood for continuous-valued variables. *IEEE Transactions on Neural Networks*, **18**, 1529–1531. [621](#)
- Hyvärinen, A. (2007b). Some extensions of score matching. *Computational Statistics and Data Analysis*, **51**, 2499–2512. [621](#)
- Hyvärinen, A. and Hoyer, P. O. (1999). Emergence of topography and complex cell properties from natural images using extensions of ica. In *NIPS*, pages 827–833. [496](#)
- Hyvärinen, A. and Pajunen, P. (1999). Nonlinear independent component analysis: Existence and uniqueness results. *Neural Networks*, **12**(3), 429–439. [496](#)
- Hyvärinen, A., Karhunen, J., and Oja, E. (2001a). *Independent Component Analysis*. Wiley-Interscience. [494](#)
- Hyvärinen, A., Hoyer, P. O., and Inki, M. O. (2001b). Topographic independent component analysis. *Neural Computation*, **13**(7), 1527–1558. [496](#)
- Hyvärinen, A., Hurri, J., and Hoyer, P. O. (2009). *Natural Image Statistics: A probabilistic approach to early computational vision*. Springer-Verlag. [371](#)

- Iba, Y. (2001). Extended ensemble Monte Carlo. *International Journal of Modern Physics*, **C12**, 623–656. [606](#)
- Inayoshi, H. and Kurita, T. (2005). Improved generalization by adding both auto-association and hidden-layer noise to neural-network-based-classifiers. *IEEE Workshop on Machine Learning for Signal Processing*, pages 141—146. [518](#)
- Ioffe, S. and Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. [100](#), [318](#), [321](#)
- Jacobs, R. A. (1988). Increased rates of convergence through learning rate adaptation. *Neural networks*, **1**(4), 295–307. [307](#)
- Jacobs, R. A., Jordan, M. I., Nowlan, S. J., and Hinton, G. E. (1991). Adaptive mixtures of local experts. *Neural Computation*, **3**, 79–87. [188](#), [453](#)
- Jaeger, H. (2003). Adaptive nonlinear system identification with echo state networks. In *Advances in Neural Information Processing Systems 15*. [406](#)
- Jaeger, H. (2007a). Discovering multiscale dynamical features with hierarchical echo state networks. Technical report, Jacobs University. [401](#)
- Jaeger, H. (2007b). Echo state network. *Scholarpedia*, **2**(9), 2330. [406](#)
- Jaeger, H. (2012). Long short-term memory in echo state networks: Details of a simulation study. Technical report, Technical report, Jacobs University Bremen. [407](#)
- Jaeger, H. and Haas, H. (2004). Harnessing nonlinearity: Predicting chaotic systems and saving energy in wireless communication. *Science*, **304**(5667), 78–80. [27](#), [406](#)
- Jaeger, H., Lukosevicius, M., Popovici, D., and Siewert, U. (2007). Optimization and applications of echo state networks with leaky- integrator neurons. *Neural Networks*, **20**(3), 335–352. [410](#)
- Jain, V., Murray, J. F., Roth, F., Turaga, S., Zhigulin, V., Briggman, K. L., Helmstaedter, M. N., Denk, W., and Seung, H. S. (2007). Supervised learning of image restoration with convolutional networks. In *Computer Vision, 2007. ICCV 2007. IEEE 11th International Conference on*, pages 1–8. IEEE. [360](#)
- Jaitly, N. and Hinton, G. (2011). Learning a better representation of speech soundwaves using restricted Boltzmann machines. In *Acoustics, Speech and Signal Processing (ICASSP), 2011 IEEE International Conference on*, pages 5884–5887. IEEE. [461](#)
- Jaitly, N. and Hinton, G. E. (2013). Vocal tract length perturbation (VTLP) improves speech recognition. In *ICML'2013*. [241](#)
- Jarrett, K., Kavukcuoglu, K., Ranzato, M., and LeCun, Y. (2009). What is the best multi-stage architecture for object recognition? In *ICCV'09*. [16](#), [24](#), [27](#), [173](#), [192](#), [226](#), [364](#), [365](#), [526](#)

- Jarzynski, C. (1997). Nonequilibrium equality for free energy differences. *Phys. Rev. Lett.*, **78**, 2690–2693. [628](#), [631](#)
- Jaynes, E. T. (2003). *Probability Theory: The Logic of Science*. Cambridge University Press. [53](#)
- Jean, S., Cho, K., Memisevic, R., and Bengio, Y. (2014). On using very large target vocabulary for neural machine translation. *arXiv:1412.2007*. [477](#)
- Jelinek, F. and Mercer, R. L. (1980). Interpolated estimation of Markov source parameters from sparse data. In E. S. Gelsema and L. N. Kanal, editors, *Pattern Recognition in Practice*. North-Holland, Amsterdam. [465](#), [476](#)
- Jia, Y. (2013). Caffe: An open source convolutional architecture for fast feature embedding. <http://caffe.berkeleyvision.org/>. [25](#), [210](#)
- Jia, Y., Huang, C., and Darrell, T. (2012). Beyond spatial pyramids: Receptive field learning for pooled image features. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 3370–3377. IEEE. [346](#)
- Jim, K.-C., Giles, C. L., and Horne, B. G. (1996). An analysis of noise in recurrent neural networks: convergence and generalization. *IEEE Transactions on Neural Networks*, **7**(6), 1424–1438. [242](#)
- Jordan, M. I. (1998). *Learning in Graphical Models*. Kluwer, Dordrecht, Netherlands. [18](#)
- Joulin, A. and Mikolov, T. (2015). Inferring algorithmic patterns with stack-augmented recurrent nets. *arXiv preprint arXiv:1503.01007*. [421](#)
- Jozefowicz, R., Zaremba, W., and Sutskever, I. (2015). An empirical evaluation of recurrent network architectures. In *ICML'2015*. [306](#), [414](#), [415](#)
- Judd, J. S. (1989). *Neural Network Design and the Complexity of Learning*. MIT press. [293](#)
- Jutten, C. and Herault, J. (1991). Blind separation of sources, part I: an adaptive algorithm based on neuromimetic architecture. *Signal Processing*, **24**, 1–10. [494](#)
- Kahou, S. E., Pal, C., Bouthillier, X., Froumenty, P., Gülcühre, c., Memisevic, R., Vincent, P., Courville, A., Bengio, Y., Ferrari, R. C., Mirza, M., Jean, S., Carrier, P. L., Dauphin, Y., Boulanger-Lewandowski, N., Aggarwal, A., Zumer, J., Lamblin, P., Raymond, J.-P., Desjardins, G., Pascanu, R., Warde-Farley, D., Torabi, A., Sharma, A., Bengio, E., Côté, M., Konda, K. R., and Wu, Z. (2013). Combining modality specific deep neural networks for emotion recognition in video. In *Proceedings of the 15th ACM on International Conference on Multimodal Interaction*. [200](#)
- Kalchbrenner, N. and Blunsom, P. (2013). Recurrent continuous translation models. In *EMNLP'2013*. [477](#)

- Kalchbrenner, N., Danihelka, I., and Graves, A. (2015). Grid long short-term memory. *arXiv preprint arXiv:1507.01526*. 397
- Kamyshanska, H. and Memisevic, R. (2015). The potential energy of an autoencoder. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. 518
- Karpathy, A. and Li, F.-F. (2015). Deep visual-semantic alignments for generating image descriptions. In *CVPR'2015*. arXiv:1412.2306. 102
- Karpathy, A., Toderici, G., Shetty, S., Leung, T., Sukthankar, R., and Fei-Fei, L. (2014). Large-scale video classification with convolutional neural networks. In *CVPR*. 21
- Karush, W. (1939). *Minima of Functions of Several Variables with Inequalities as Side Constraints*. Master's thesis, Dept. of Mathematics, Univ. of Chicago. 95
- Katz, S. M. (1987). Estimation of probabilities from sparse data for the language model component of a speech recognizer. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, **ASSP-35**(3), 400–401. 465, 476
- Kavukcuoglu, K., Ranzato, M., and LeCun, Y. (2008). Fast inference in sparse coding algorithms with applications to object recognition. Technical report, Computational and Biological Learning Lab, Courant Institute, NYU. Tech Report CBLL-TR-2008-12-01. 526
- Kavukcuoglu, K., Ranzato, M.-A., Fergus, R., and LeCun, Y. (2009). Learning invariant features through topographic filter maps. In *CVPR'2009*. 526
- Kavukcuoglu, K., Sermanet, P., Boureau, Y.-L., Gregor, K., Mathieu, M., and LeCun, Y. (2010). Learning convolutional feature hierarchies for visual recognition. In *NIPS'2010*. 365, 526
- Kelley, H. J. (1960). Gradient theory of optimal flight paths. *ARS Journal*, **30**(10), 947–954. 225
- Khan, F., Zhu, X., and Mutlu, B. (2011). How do humans teach: On curriculum learning and teaching dimension. In *Advances in Neural Information Processing Systems 24 (NIPS'11)*, pages 1449–1457. 329
- Kim, S. K., McAfee, L. C., McMahon, P. L., and Olukotun, K. (2009). A highly scalable restricted Boltzmann machine FPGA implementation. In *Field Programmable Logic and Applications, 2009. FPL 2009. International Conference on*, pages 367–372. IEEE. 454
- Kindermann, R. (1980). *Markov Random Fields and Their Applications (Contemporary Mathematics ; V. 1)*. American Mathematical Society. 569
- Kingma, D. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*. 308

- Kingma, D. and LeCun, Y. (2010). Regularized estimation of image statistics by score matching. In *NIPS'2010*. 516, 623
- Kingma, D., Rezende, D., Mohamed, S., and Welling, M. (2014). Semi-supervised learning with deep generative models. In *NIPS'2014*. 429
- Kingma, D. P. (2013). Fast gradient-based inference with continuous latent variable models in auxiliary form. Technical report, arxiv:1306.0733. 655, 691, 698
- Kingma, D. P. and Welling, M. (2014a). Auto-encoding variational bayes. In *Proceedings of the International Conference on Learning Representations (ICLR)*. 691, 701
- Kingma, D. P. and Welling, M. (2014b). Efficient gradient-based inference through transformations between bayes nets and neural nets. Technical report, arxiv:1402.0480. 691
- Kirkpatrick, S., Jr., C. D. G., , and Vecchi, M. P. (1983). Optimization by simulated annealing. *Science*, 220, 671–680. 328
- Kiros, R., Salakhutdinov, R., and Zemel, R. (2014a). Multimodal neural language models. In *ICML'2014*. 102
- Kiros, R., Salakhutdinov, R., and Zemel, R. (2014b). Unifying visual-semantic embeddings with multimodal neural language models. *arXiv:1411.2539 [cs.LG]*. 102, 411
- Klementiev, A., Titov, I., and Bhattacharyya, B. (2012). Inducing crosslingual distributed representations of words. In *Proceedings of COLING 2012*. 479, 542
- Knowles-Barley, S., Jones, T. R., Morgan, J., Lee, D., Kasthuri, N., Lichtman, J. W., and Pfister, H. (2014). Deep learning for the connectome. *GPU Technology Conference*. 26
- Koller, D. and Friedman, N. (2009). *Probabilistic Graphical Models: Principles and Techniques*. MIT Press. 585, 598, 648
- Konig, Y., Bourlard, H., and Morgan, N. (1996). REMAP: Recursive estimation and maximization of a posteriori probabilities – application to transition-based connectionist speech recognition. In D. Touretzky, M. Mozer, and M. Hasselmo, editors, *Advances in Neural Information Processing Systems 8 (NIPS'95)*. MIT Press, Cambridge, MA. 462
- Koren, Y. (2009). The BellKor solution to the Netflix grand prize. 256, 482
- Kotzias, D., Denil, M., de Freitas, N., and Smyth, P. (2015). From group to individual labels using deep features. In *ACM SIGKDD*. 106
- Koutnik, J., Greff, K., Gomez, F., and Schmidhuber, J. (2014). A clockwork RNN. In *ICML'2014*. 410
- Kočiský, T., Hermann, K. M., and Blunsom, P. (2014). Learning Bilingual Word Representations by Marginalizing Alignments. In *Proceedings of ACL*. 479

- Krause, O., Fischer, A., Glasmachers, T., and Igel, C. (2013). Approximation properties of DBNs with binary hidden units and real-valued visible units. In *ICML'2013*. 556
- Krizhevsky, A. (2010). Convolutional deep belief networks on CIFAR-10. Technical report, University of Toronto. Unpublished Manuscript: <http://www.cs.utoronto.ca/~kriz/conv-cifar10-aug2010.pdf>. 449
- Krizhevsky, A. and Hinton, G. (2009). Learning multiple layers of features from tiny images. Technical report, University of Toronto. 21, 564
- Krizhevsky, A. and Hinton, G. E. (2011). Using very deep autoencoders for content-based image retrieval. In *ESANN*. 528
- Krizhevsky, A., Sutskever, I., and Hinton, G. (2012). ImageNet classification with deep convolutional neural networks. In *NIPS'2012*. 23, 24, 27, 100, 200, 372, 457, 461
- Krueger, K. A. and Dayan, P. (2009). Flexible shaping: how learning in small steps helps. *Cognition*, 110, 380–394. 329
- Kuhn, H. W. and Tucker, A. W. (1951). Nonlinear programming. In *Proceedings of the Second Berkeley Symposium on Mathematical Statistics and Probability*, pages 481–492, Berkeley, Calif. University of California Press. 95
- Kumar, A., Irsay, O., Su, J., Bradbury, J., English, R., Pierce, B., Ondruska, P., Iyyer, M., Gulrajani, I., and Socher, R. (2015). Ask me anything: Dynamic memory networks for natural language processing. *arXiv:1506.07285*. 421, 488
- Kumar, M. P., Packer, B., and Koller, D. (2010). Self-paced learning for latent variable models. In *NIPS'2010*. 329
- Lang, K. J. and Hinton, G. E. (1988). The development of the time-delay neural network architecture for speech recognition. Technical Report CMU-CS-88-152, Carnegie-Mellon University. 368, 375, 409
- Lang, K. J., Waibel, A. H., and Hinton, G. E. (1990). A time-delay neural network architecture for isolated word recognition. *Neural networks*, 3(1), 23–43. 375
- Langford, J. and Zhang, T. (2008). The epoch-greedy algorithm for contextual multi-armed bandits. In *NIPS'2008*, pages 1096—1103. 483
- Lappalainen, H., Giannakopoulos, X., Honkela, A., and Karhunen, J. (2000). Nonlinear independent component analysis using ensemble learning: Experiments and discussion. In *Proc. ICA*. Citeseer. 496
- Larochelle, H. and Bengio, Y. (2008). Classification using discriminative restricted Boltzmann machines. In *ICML'2008*. 244, 254, 533, 688, 717

- Larochelle, H. and Hinton, G. E. (2010). Learning to combine foveal glimpses with a third-order Boltzmann machine. In *Advances in Neural Information Processing Systems 23*, pages 1243–1251. [368](#)
- Larochelle, H. and Murray, I. (2011). The Neural Autoregressive Distribution Estimator. In *AISTATS'2011*. [707](#), [710](#)
- Larochelle, H., Erhan, D., and Bengio, Y. (2008). Zero-data learning of new tasks. In *AAAI Conference on Artificial Intelligence*. [542](#)
- Larochelle, H., Bengio, Y., Louradour, J., and Lamblin, P. (2009). Exploring strategies for training deep neural networks. *Journal of Machine Learning Research*, **10**, 1–40. [538](#)
- Lasserre, J. A., Bishop, C. M., and Minka, T. P. (2006). Principled hybrids of generative and discriminative models. In *Proceedings of the Computer Vision and Pattern Recognition Conference (CVPR'06)*, pages 87–94, Washington, DC, USA. IEEE Computer Society. [244](#), [252](#)
- Le, Q., Ngiam, J., Chen, Z., hao Chia, D. J., Koh, P. W., and Ng, A. (2010). Tiled convolutional neural networks. In J. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. Zemel, and A. Culotta, editors, *Advances in Neural Information Processing Systems 23 (NIPS'10)*, pages 1279–1287. [353](#)
- Le, Q., Ngiam, J., Coates, A., Lahiri, A., Prochnow, B., and Ng, A. (2011). On optimization methods for deep learning. In *Proc. ICML'2011*. ACM. [316](#)
- Le, Q., Ranzato, M., Monga, R., Devin, M., Corrado, G., Chen, K., Dean, J., and Ng, A. (2012). Building high-level features using large scale unsupervised learning. In *ICML'2012*. [24](#), [27](#)
- Le Roux, N. and Bengio, Y. (2008). Representational power of restricted Boltzmann machines and deep belief networks. *Neural Computation*, **20**(6), 1631–1649. [556](#), [657](#)
- Le Roux, N. and Bengio, Y. (2010). Deep belief networks are compact universal approximators. *Neural Computation*, **22**(8), 2192–2207. [556](#)
- LeCun, Y. (1985). Une procédure d'apprentissage pour Réseau à seuil assymétrique. In *Cognitiva 85: A la Frontière de l'Intelligence Artificielle, des Sciences de la Connaissance et des Neurosciences*, pages 599–604, Paris 1985. CESTA, Paris. [225](#)
- LeCun, Y. (1986). Learning processes in an asymmetric threshold network. In F. Fogelman-Soulie, E. Bienenstock, and G. Weisbuch, editors, *Disordered Systems and Biological Organization*, pages 233–240. Springer-Verlag, Les Houches, France. [351](#)
- LeCun, Y. (1987). *Modèles connexionnistes de l'apprentissage*. Ph.D. thesis, Université de Paris VI. [18](#), [505](#), [518](#)
- LeCun, Y. (1989). Generalization and network design strategies. Technical Report CRG-TR-89-4, University of Toronto. [331](#), [351](#)

- LeCun, Y., Jackel, L. D., Boser, B., Denker, J. S., Graf, H. P., Guyon, I., Henderson, D., Howard, R. E., and Hubbard, W. (1989). Handwritten digit recognition: Applications of neural network chips and automatic learning. *IEEE Communications Magazine*, **27**(11), 41–46. [369](#)
- LeCun, Y., Bottou, L., Orr, G. B., and Müller, K.-R. (1998a). Efficient backprop. In *Neural Networks, Tricks of the Trade*, Lecture Notes in Computer Science LNCS 1524. Springer Verlag. [310](#), [432](#)
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998b). Gradient based learning applied to document recognition. *Proc. IEEE*. [16](#), [18](#), [21](#), [27](#), [372](#), [461](#), [463](#)
- LeCun, Y., Kavukcuoglu, K., and Farabet, C. (2010). Convolutional networks and applications in vision. In *Circuits and Systems (ISCAS), Proceedings of 2010 IEEE International Symposium on*, pages 253–256. IEEE. [372](#)
- L’Ecuyer, P. (1994). Efficiency improvement and variance reduction. In *Proceedings of the 1994 Winter Simulation Conference*, pages 122—132. [692](#)
- Lee, C.-Y., Xie, S., Gallagher, P., Zhang, Z., and Tu, Z. (2014). Deeply-supervised nets. *arXiv preprint arXiv:1409.5185*. [327](#)
- Lee, H., Battle, A., Raina, R., and Ng, A. (2007). Efficient sparse coding algorithms. In B. Schölkopf, J. Platt, and T. Hoffman, editors, *Advances in Neural Information Processing Systems 19 (NIPS’06)*, pages 801–808. MIT Press. [640](#)
- Lee, H., Ekanadham, C., and Ng, A. (2008). Sparse deep belief net model for visual area V2. In *NIPS’07*. [254](#)
- Lee, H., Grosse, R., Ranganath, R., and Ng, A. Y. (2009). Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. In L. Bottou and M. Littman, editors, *Proceedings of the Twenty-sixth International Conference on Machine Learning (ICML’09)*. ACM, Montreal, Canada. [364](#), [685](#), [686](#)
- Lee, Y. J. and Grauman, K. (2011). Learning the easy things first: self-paced visual category discovery. In *CVPR’2011*. [329](#)
- Leibniz, G. W. (1676). Memoir using the chain rule. (Cited in TMME 7:2&3 p 321-332, 2010). [224](#)
- Lenat, D. B. and Guha, R. V. (1989). *Building large knowledge-based systems; representation and inference in the Cyc project*. Addison-Wesley Longman Publishing Co., Inc. [2](#)
- Leshno, M., Lin, V. Y., Pinkus, A., and Schocken, S. (1993). Multilayer feedforward networks with a nonpolynomial activation function can approximate any function. *Neural Networks*, **6**, 861—867. [197](#), [198](#)

- Levenberg, K. (1944). A method for the solution of certain non-linear problems in least squares. *Quarterly Journal of Applied Mathematics*, **II**(2), 164–168. [312](#)
- L'Hôpital, G. F. A. (1696). *Analyse des infiniment petits, pour l'intelligence des lignes courbes*. Paris: L'Imprimerie Royale. [224](#)
- Li, Y., Swersky, K., and Zemel, R. S. (2015). Generative moment matching networks. *CoRR*, **abs/1502.02761**. [705](#)
- Lin, T., Horne, B. G., Tino, P., and Giles, C. L. (1996). Learning long-term dependencies is not as difficult with NARX recurrent neural networks. *IEEE Transactions on Neural Networks*, **7**(6), 1329–1338. [409](#)
- Lin, Y., Liu, Z., Sun, M., Liu, Y., and Zhu, X. (2015). Learning entity and relation embeddings for knowledge graph completion. In *Proc. AAAI'15*. [487](#)
- Linde, N. (1992). The machine that changed the world, episode 3. Documentary miniseries. [2](#)
- Lindsey, C. and Lindblad, T. (1994). Review of hardware neural networks: a user's perspective. In *Proc. Third Workshop on Neural Networks: From Biology to High Energy Physics*, pages 195—202, Isola d'Elba, Italy. [454](#)
- Linnainmaa, S. (1976). Taylor expansion of the accumulated rounding error. *BIT Numerical Mathematics*, **16**(2), 146–160. [225](#)
- LISA (2008). Deep learning tutorials: Restricted Boltzmann machines. Technical report, LISA Lab, Université de Montréal. [591](#)
- Long, P. M. and Servedio, R. A. (2010). Restricted Boltzmann machines are hard to approximately evaluate or simulate. In *Proceedings of the 27th International Conference on Machine Learning (ICML'10)*. [660](#)
- Lotter, W., Kreiman, G., and Cox, D. (2015). Unsupervised learning of visual structure using predictive generative networks. *arXiv preprint arXiv:1511.06380*. [547](#), [548](#)
- Lovelace, A. (1842). Notes upon L. F. Menabrea's “Sketch of the Analytical Engine invented by Charles Babbage”. [1](#)
- Lu, L., Zhang, X., Cho, K., and Renals, S. (2015). A study of the recurrent neural network encoder-decoder for large vocabulary speech recognition. In *Proc. Interspeech*. [463](#)
- Lu, T., Pál, D., and Pál, M. (2010). Contextual multi-armed bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 485–492. [483](#)
- Luenberger, D. G. (1984). *Linear and Nonlinear Programming*. Addison Wesley. [317](#)
- Lukoševičius, M. and Jaeger, H. (2009). Reservoir computing approaches to recurrent neural network training. *Computer Science Review*, **3**(3), 127–149. [406](#)

- Luo, H., Shen, R., Niu, C., and Ullrich, C. (2011). Learning class-relevant features and class-irrelevant features via a hybrid third-order RBM. In *International Conference on Artificial Intelligence and Statistics*, pages 470–478. **689**
- Luo, H., Carrier, P. L., Courville, A., and Bengio, Y. (2013). Texture modeling with convolutional spike-and-slab RBMs and deep extensions. In *AISTATS’2013*. **102**
- Lyu, S. (2009). Interpretation and generalization of score matching. In *Proceedings of the Twenty-fifth Conference in Uncertainty in Artificial Intelligence (UAI’09)*. **621**
- Ma, J., Sheridan, R. P., Liaw, A., Dahl, G. E., and Svetnik, V. (2015). Deep neural nets as a method for quantitative structure – activity relationships. *J. Chemical information and modeling*. **533**
- Maas, A. L., Hannun, A. Y., and Ng, A. Y. (2013). Rectifier nonlinearities improve neural network acoustic models. In *ICML Workshop on Deep Learning for Audio, Speech, and Language Processing*. **192**
- Maass, W. (1992). Bounds for the computational power and learning complexity of analog neural nets (extended abstract). In *Proc. of the 25th ACM Symp. Theory of Computing*, pages 335–344. **198**
- Maass, W., Schnitger, G., and Sontag, E. D. (1994). A comparison of the computational power of sigmoid and Boolean threshold circuits. *Theoretical Advances in Neural Computation and Learning*, pages 127–151. **198**
- Maass, W., Natschlaeger, T., and Markram, H. (2002). Real-time computing without stable states: A new framework for neural computation based on perturbations. *Neural Computation*, **14**(11), 2531–2560. **406**
- MacKay, D. (2003). *Information Theory, Inference and Learning Algorithms*. Cambridge University Press. **73**
- Maclaurin, D., Duvenaud, D., and Adams, R. P. (2015). Gradient-based hyperparameter optimization through reversible learning. *arXiv preprint arXiv:1502.03492*. **438**
- Mao, J., Xu, W., Yang, Y., Wang, J., Huang, Z., and Yuille, A. L. (2015). Deep captioning with multimodal recurrent neural networks. In *ICLR’2015*. arXiv:1410.1090. **102**
- Marcotte, P. and Savard, G. (1992). Novel approaches to the discrimination problem. *Zeitschrift für Operations Research (Theory)*, **36**, 517–545. **276**
- Marlin, B. and de Freitas, N. (2011). Asymptotic efficiency of deterministic estimators for discrete energy-based models: Ratio matching and pseudolikelihood. In *UAI’2011*. **620**, **622**

- Marlin, B., Swersky, K., Chen, B., and de Freitas, N. (2010). Inductive principles for restricted Boltzmann machine learning. In *Proceedings of The Thirteenth International Conference on Artificial Intelligence and Statistics (AISTATS'10)*, volume 9, pages 509–516. [616](#), [621](#), [622](#)
- Marquardt, D. W. (1963). An algorithm for least-squares estimation of non-linear parameters. *Journal of the Society of Industrial and Applied Mathematics*, **11**(2), 431–441. [312](#)
- Marr, D. and Poggio, T. (1976). Cooperative computation of stereo disparity. *Science*, **194**. [368](#)
- Martens, J. (2010). Deep learning via Hessian-free optimization. In L. Bottou and M. Littman, editors, *Proceedings of the Twenty-seventh International Conference on Machine Learning (ICML-10)*, pages 735–742. ACM. [304](#)
- Martens, J. and Medabalimi, V. (2014). On the expressive efficiency of sum product networks. *arXiv:1411.7717*. [557](#)
- Martens, J. and Sutskever, I. (2011). Learning recurrent neural networks with Hessian-free optimization. In *Proc. ICML'2011*. ACM. [415](#)
- Mase, S. (1995). Consistency of the maximum pseudo-likelihood estimator of continuous state space Gibbsian processes. *The Annals of Applied Probability*, **5**(3), pp. 603–612. [619](#)
- McClelland, J., Rumelhart, D., and Hinton, G. (1995). The appeal of parallel distributed processing. In *Computation & intelligence*, pages 305–341. American Association for Artificial Intelligence. [17](#)
- McCulloch, W. S. and Pitts, W. (1943). A logical calculus of ideas immanent in nervous activity. *Bulletin of Mathematical Biophysics*, **5**, 115–133. [14](#), [15](#)
- Mead, C. and Ismail, M. (2012). *Analog VLSI implementation of neural systems*, volume 80. Springer Science & Business Media. [454](#)
- Melchior, J., Fischer, A., and Wiskott, L. (2013). How to center binary deep Boltzmann machines. *arXiv preprint arXiv:1311.1354*. [675](#)
- Memisevic, R. and Hinton, G. E. (2007). Unsupervised learning of image transformations. In *Proceedings of the Computer Vision and Pattern Recognition Conference (CVPR'07)*. [688](#)
- Memisevic, R. and Hinton, G. E. (2010). Learning to represent spatial transformations with factored higher-order Boltzmann machines. *Neural Computation*, **22**(6), 1473–1492. [688](#)

- Mesnil, G., Dauphin, Y., Glorot, X., Rifai, S., Bengio, Y., Goodfellow, I., Lavoie, E., Muller, X., Desjardins, G., Warde-Farley, D., Vincent, P., Courville, A., and Bergstra, J. (2011). Unsupervised and transfer learning challenge: a deep learning approach. In *JMLR W&CP: Proc. Unsupervised and Transfer Learning*, volume 7. **200**, 535, 541
- Mesnil, G., Rifai, S., Dauphin, Y., Bengio, Y., and Vincent, P. (2012). Surfing on the manifold. Learning Workshop, Snowbird. **713**
- Miikkulainen, R. and Dyer, M. G. (1991). Natural language processing with modular PDP networks and distributed lexicon. *Cognitive Science*, **15**, 343–399. **480**
- Mikolov, T. (2012). *Statistical Language Models based on Neural Networks*. Ph.D. thesis, Brno University of Technology. **417**
- Mikolov, T., Deoras, A., Kombrink, S., Burget, L., and Cernocky, J. (2011a). Empirical evaluation and combination of advanced language modeling techniques. In *Proc. 12th annual conference of the international speech communication association (INTERSPEECH 2011)*. **475**
- Mikolov, T., Deoras, A., Povey, D., Burget, L., and Cernocky, J. (2011b). Strategies for training large scale neural network language models. In *Proc. ASRU'2011*. **329**, 475
- Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013a). Efficient estimation of word representations in vector space. In *International Conference on Learning Representations: Workshops Track*. **539**
- Mikolov, T., Le, Q. V., and Sutskever, I. (2013b). Exploiting similarities among languages for machine translation. Technical report, arXiv:1309.4168. **542**
- Minka, T. (2005). Divergence measures and message passing. *Microsoft Research Cambridge UK Tech Rep MSRTR2005173*, **72**(TR-2005-173). **628**
- Minsky, M. L. and Papert, S. A. (1969). *Perceptrons*. MIT Press, Cambridge. **15**
- Mirza, M. and Osindero, S. (2014). Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*. **703**
- Mishkin, D. and Matas, J. (2015). All you need is a good init. *arXiv preprint arXiv:1511.06422*. **305**
- Misra, J. and Saha, I. (2010). Artificial neural networks in hardware: A survey of two decades of progress. *Neurocomputing*, **74**(1), 239–255. **454**
- Mitchell, T. M. (1997). *Machine Learning*. McGraw-Hill, New York. **99**
- Miyato, T., Maeda, S., Koyama, M., Nakae, K., and Ishii, S. (2015). Distributional smoothing with virtual adversarial training. In *ICLR*. Preprint: arXiv:1507.00677. **268**

- Mnih, A. and Gregor, K. (2014). Neural variational inference and learning in belief networks. In *ICML'2014*. 693, 695
- Mnih, A. and Hinton, G. E. (2007). Three new graphical models for statistical language modelling. In Z. Ghahramani, editor, *Proceedings of the Twenty-fourth International Conference on Machine Learning (ICML'07)*, pages 641–648. ACM. 467
- Mnih, A. and Hinton, G. E. (2009). A scalable hierarchical distributed language model. In D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, editors, *Advances in Neural Information Processing Systems 21 (NIPS'08)*, pages 1081–1088. 470
- Mnih, A. and Kavukcuoglu, K. (2013). Learning word embeddings efficiently with noise-contrastive estimation. In C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 2265–2273. Curran Associates, Inc. 475, 625
- Mnih, A. and Teh, Y. W. (2012). A fast and simple algorithm for training neural probabilistic language models. In *ICML'2012*, pages 1751–1758. 475
- Mnih, V. and Hinton, G. (2010). Learning to detect roads in high-resolution aerial images. In *Proceedings of the 11th European Conference on Computer Vision (ECCV)*. 102
- Mnih, V., Larochelle, H., and Hinton, G. (2011). Conditional restricted Boltzmann machines for structure output prediction. In *Proc. Conf. on Uncertainty in Artificial Intelligence (UAI)*. 687
- Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., and Wierstra, D. (2013). Playing Atari with deep reinforcement learning. Technical report, arXiv:1312.5602. 106
- Mnih, V., Heess, N., Graves, A., and Kavukcuoglu, K. (2014). Recurrent models of visual attention. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Weinberger, editors, *NIPS'2014*, pages 2204–2212. 693
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidgeland, A. K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., and Hassabis, D. (2015). Human-level control through deep reinforcement learning. *Nature*, 518, 529–533. 25
- Mobahi, H. and Fisher, III, J. W. (2015). A theoretical analysis of optimization by Gaussian continuation. In *AAAI'2015*. 328
- Mobahi, H., Collobert, R., and Weston, J. (2009). Deep learning from temporal coherence in video. In L. Bottou and M. Littman, editors, *Proceedings of the 26th International Conference on Machine Learning*, pages 737–744, Montreal. Omnipress. 497
- Mohamed, A., Dahl, G., and Hinton, G. (2009). Deep belief networks for phone recognition. 462

- Mohamed, A., Sainath, T. N., Dahl, G., Ramabhadran, B., Hinton, G. E., and Picheny, M. A. (2011). Deep belief networks using discriminative features for phone recognition. In *Acoustics, Speech and Signal Processing (ICASSP), 2011 IEEE International Conference on*, pages 5060–5063. IEEE. [462](#)
- Mohamed, A., Dahl, G., and Hinton, G. (2012a). Acoustic modeling using deep belief networks. *IEEE Trans. on Audio, Speech and Language Processing*, **20**(1), 14–22. [462](#)
- Mohamed, A., Hinton, G., and Penn, G. (2012b). Understanding how deep belief networks perform acoustic modelling. In *Acoustics, Speech and Signal Processing (ICASSP), 2012 IEEE International Conference on*, pages 4273–4276. IEEE. [462](#)
- Moller, M. F. (1993). A scaled conjugate gradient algorithm for fast supervised learning. *Neural Networks*, **6**, 525–533. [316](#)
- Montavon, G. and Muller, K.-R. (2012). Deep Boltzmann machines and the centering trick. In G. Montavon, G. Orr, and K.-R. Müller, editors, *Neural Networks: Tricks of the Trade*, volume 7700 of *Lecture Notes in Computer Science*, pages 621–637. Preprint: <http://arxiv.org/abs/1203.3783>. [675](#)
- Montúfar, G. (2014). Universal approximation depth and errors of narrow belief networks with discrete units. *Neural Computation*, **26**. [556](#)
- Montúfar, G. and Ay, N. (2011). Refinements of universal approximation results for deep belief networks and restricted Boltzmann machines. *Neural Computation*, **23**(5), 1306–1319. [556](#)
- Montufar, G. F., Pascanu, R., Cho, K., and Bengio, Y. (2014). On the number of linear regions of deep neural networks. In *NIPS’2014*. [19](#), [199](#)
- Mor-Yosef, S., Samueloff, A., Modan, B., Navot, D., and Schenker, J. G. (1990). Ranking the risk factors for cesarean: logistic regression analysis of a nationwide study. *Obstet Gynecol*, **75**(6), 944–7. [3](#)
- Morin, F. and Bengio, Y. (2005). Hierarchical probabilistic neural network language model. In *AISTATS’2005*. [470](#), [472](#)
- Mozer, M. C. (1992). The induction of multiscale temporal structure. In J. M. S. Hanson and R. Lippmann, editors, *Advances in Neural Information Processing Systems 4 (NIPS’91)*, pages 275–282, San Mateo, CA. Morgan Kaufmann. [410](#)
- Murphy, K. P. (2012). *Machine Learning: a Probabilistic Perspective*. MIT Press, Cambridge, MA, USA. [62](#), [98](#), [145](#)
- Murray, B. U. I. and Larochelle, H. (2014). A deep and tractable density estimator. In *ICML’2014*. [189](#), [712](#)
- Nair, V. and Hinton, G. (2010). Rectified linear units improve restricted Boltzmann machines. In *ICML’2010*. [16](#), [173](#), [196](#)

- Nair, V. and Hinton, G. E. (2009). 3d object recognition with deep belief nets. In Y. Bengio, D. Schuurmans, J. D. Lafferty, C. K. I. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 1339–1347. Curran Associates, Inc. [688](#)
- Narayanan, H. and Mitter, S. (2010). Sample complexity of testing the manifold hypothesis. In *NIPS’2010*. [163](#)
- Naumann, U. (2008). Optimal Jacobian accumulation is NP-complete. *Mathematical Programming*, **112**(2), 427–441. [221](#)
- Navigli, R. and Velardi, P. (2005). Structural semantic interconnections: a knowledge-based approach to word sense disambiguation. *IEEE Trans. Pattern Analysis and Machine Intelligence*, **27**(7), 1075–1086. [487](#)
- Neal, R. and Hinton, G. (1999). A view of the EM algorithm that justifies incremental, sparse, and other variants. In M. I. Jordan, editor, *Learning in Graphical Models*. MIT Press, Cambridge, MA. [637](#)
- Neal, R. M. (1990). Learning stochastic feedforward networks. Technical report. [694](#)
- Neal, R. M. (1993). Probabilistic inference using Markov chain Monte-Carlo methods. Technical Report CRG-TR-93-1, Dept. of Computer Science, University of Toronto. [682](#)
- Neal, R. M. (1994). Sampling from multimodal distributions using tempered transitions. Technical Report 9421, Dept. of Statistics, University of Toronto. [606](#)
- Neal, R. M. (1996). *Bayesian Learning for Neural Networks*. Lecture Notes in Statistics. Springer. [264](#)
- Neal, R. M. (2001). Annealed importance sampling. *Statistics and Computing*, **11**(2), 125–139. [628](#), [630](#), [631](#), [632](#)
- Neal, R. M. (2005). Estimating ratios of normalizing constants using linked importance sampling. [632](#)
- Nesterov, Y. (1983). A method of solving a convex programming problem with convergence rate  $O(1/k^2)$ . *Soviet Mathematics Doklady*, **27**, 372–376. [300](#)
- Nesterov, Y. (2004). *Introductory lectures on convex optimization : a basic course*. Applied optimization. Kluwer Academic Publ., Boston, Dordrecht, London. [300](#)
- Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., and Ng, A. Y. (2011). Reading digits in natural images with unsupervised feature learning. Deep Learning and Unsupervised Feature Learning Workshop, NIPS. [21](#)
- Ney, H. and Kneser, R. (1993). Improved clustering techniques for class-based statistical language modelling. In *European Conference on Speech Communication and Technology (Eurospeech)*, pages 973–976, Berlin. [466](#)

- Ng, A. (2015). Advice for applying machine learning. <https://see.stanford.edu/materials/aimlcs229/ML-advice.pdf>. 424
- Niesler, T. R., Whittaker, E. W. D., and Woodland, P. C. (1998). Comparison of part-of-speech and automatically derived category-based language models for speech recognition. In *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 177–180. 466
- Ning, F., Delhomme, D., LeCun, Y., Piano, F., Bottou, L., and Barbano, P. E. (2005). Toward automatic phenotyping of developing embryos from videos. *Image Processing, IEEE Transactions on*, **14**(9), 1360–1371. 361
- Nocedal, J. and Wright, S. (2006). *Numerical Optimization*. Springer. 92, 95
- Norouzi, M. and Fleet, D. J. (2011). Minimal loss hashing for compact binary codes. In *ICML'2011*. 528
- Nowlan, S. J. (1990). Competing experts: An experimental investigation of associative mixture models. Technical Report CRG-TR-90-5, University of Toronto. 453
- Nowlan, S. J. and Hinton, G. E. (1992). Simplifying neural networks by soft weight-sharing. *Neural Computation*, **4**(4), 473–493. 139
- Olshausen, B. and Field, D. J. (2005). How close are we to understanding V1? *Neural Computation*, **17**, 1665–1699. 16
- Olshausen, B. A. and Field, D. J. (1996). Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, **381**, 607–609. 146, 254, 371, 499
- Olshausen, B. A., Anderson, C. H., and Van Essen, D. C. (1993). A neurobiological model of visual attention and invariant pattern recognition based on dynamic routing of information. *J. Neurosci.*, **13**(11), 4700–4719. 453
- Opper, M. and Archambeau, C. (2009). The variational Gaussian approximation revisited. *Neural computation*, **21**(3), 786–792. 691
- Quab, M., Bottou, L., Laptev, I., and Sivic, J. (2014). Learning and transferring mid-level image representations using convolutional neural networks. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 1717–1724. IEEE. 539
- Osindero, S. and Hinton, G. E. (2008). Modeling image patches with a directed hierarchy of Markov random fields. In J. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems 20 (NIPS'07)*, pages 1121–1128, Cambridge, MA. MIT Press. 635
- Ovid and Martin, C. (2004). *Metamorphoses*. W.W. Norton. 1

- Paccanaro, A. and Hinton, G. E. (2000). Extracting distributed representations of concepts and relations from positive and negative propositions. In *International Joint Conference on Neural Networks (IJCNN)*, Como, Italy. IEEE, New York. 487
- Paine, T. L., Khorrami, P., Han, W., and Huang, T. S. (2014). An analysis of unsupervised pre-training in light of recent advances. *arXiv preprint arXiv:1412.6597*. 535
- Palatucci, M., Pomerleau, D., Hinton, G. E., and Mitchell, T. M. (2009). Zero-shot learning with semantic output codes. In Y. Bengio, D. Schuurmans, J. D. Lafferty, C. K. I. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 1410–1418. Curran Associates, Inc. 542
- Parker, D. B. (1985). Learning-logic. Technical Report TR-47, Center for Comp. Research in Economics and Management Sci., MIT. 225
- Pascanu, R., Mikolov, T., and Bengio, Y. (2013a). On the difficulty of training recurrent neural networks. In *ICML'2013*. 289, 403, 406, 410, 417, 419
- Pascanu, R., Montufar, G., and Bengio, Y. (2013b). On the number of inference regions of deep feed forward networks with piece-wise linear activations. Technical report, U. Montreal, arXiv:1312.6098. 198
- Pascanu, R., Gülcühre, Ç., Cho, K., and Bengio, Y. (2014a). How to construct deep recurrent neural networks. In *ICLR'2014*. 19, 199, 264, 399, 400, 401, 413, 463
- Pascanu, R., Montufar, G., and Bengio, Y. (2014b). On the number of inference regions of deep feed forward networks with piece-wise linear activations. In *ICLR'2014*. 553
- Pati, Y., Rezaifar, R., and Krishnaprasad, P. (1993). Orthogonal matching pursuit: Recursive function approximation with applications to wavelet decomposition. In *Proceedings of the 27 th Annual Asilomar Conference on Signals, Systems, and Computers*, pages 40–44. 254
- Pearl, J. (1985). Bayesian networks: A model of self-activated memory for evidential reasoning. In *Proceedings of the 7th Conference of the Cognitive Science Society, University of California, Irvine*, pages 329–334. 566
- Pearl, J. (1988). *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann. 54
- Perron, O. (1907). Zur theorie der matrices. *Mathematische Annalen*, **64**(2), 248–263. 600
- Petersen, K. B. and Pedersen, M. S. (2006). The matrix cookbook. Version 20051003. 31
- Peterson, G. B. (2004). A day of great illumination: B. F. Skinner's discovery of shaping. *Journal of the Experimental Analysis of Behavior*, **82**(3), 317–328. 329
- Pham, D.-T., Garat, P., and Jutten, C. (1992). Separation of a mixture of independent sources through a maximum likelihood approach. In *EUSIPCO*, pages 771–774. 494

- Pham, P.-H., Jelaca, D., Farabet, C., Martini, B., LeCun, Y., and Culurciello, E. (2012). NeuFlow: dataflow vision processing system-on-a-chip. In *Circuits and Systems (MWCAS), 2012 IEEE 55th International Midwest Symposium on*, pages 1044–1047. IEEE. [454](#)
- Pinheiro, P. H. O. and Collobert, R. (2014). Recurrent convolutional neural networks for scene labeling. In *ICML'2014*. [360](#)
- Pinheiro, P. H. O. and Collobert, R. (2015). From image-level to pixel-level labeling with convolutional networks. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. [360](#)
- Pinto, N., Cox, D. D., and DiCarlo, J. J. (2008). Why is real-world visual object recognition hard? *PLoS Comput Biol*, **4**. [459](#)
- Pinto, N., Stone, Z., Zickler, T., and Cox, D. (2011). Scaling up biologically-inspired computer vision: A case study in unconstrained face recognition on facebook. In *Computer Vision and Pattern Recognition Workshops (CVPRW), 2011 IEEE Computer Society Conference on*, pages 35–42. IEEE. [364](#)
- Pollack, J. B. (1990). Recursive distributed representations. *Artificial Intelligence*, **46**(1), 77–105. [401](#)
- Polyak, B. and Juditsky, A. (1992). Acceleration of stochastic approximation by averaging. *SIAM J. Control and Optimization*, **30**(4), 838–855. [323](#)
- Polyak, B. T. (1964). Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, **4**(5), 1–17. [296](#)
- Poole, B., Sohl-Dickstein, J., and Ganguli, S. (2014). Analyzing noise in autoencoders and deep networks. *CoRR*, **abs/1406.1831**. [241](#)
- Poon, H. and Domingos, P. (2011). Sum-product networks: A new deep architecture. In *Proceedings of the Twenty-seventh Conference in Uncertainty in Artificial Intelligence (UAI)*, Barcelona, Spain. [557](#)
- Presley, R. K. and Haggard, R. L. (1994). A fixed point implementation of the backpropagation learning algorithm. In *Southeastcon'94. Creative Technology Transfer-A Global Affair., Proceedings of the 1994 IEEE*, pages 136–138. IEEE. [454](#)
- Price, R. (1958). A useful theorem for nonlinear devices having Gaussian inputs. *IEEE Transactions on Information Theory*, **4**(2), 69–72. [691](#)
- Quiroga, R. Q., Reddy, L., Kreiman, G., Koch, C., and Fried, I. (2005). Invariant visual representation by single neurons in the human brain. *Nature*, **435**(7045), 1102–1107. [367](#)

- Radford, A., Metz, L., and Chintala, S. (2015). Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*. [555](#), [703](#), [704](#)
- Raiko, T., Yao, L., Cho, K., and Bengio, Y. (2014). Iterative neural autoregressive distribution estimator (NADE-k). Technical report, arXiv:1406.1485. [678](#), [711](#)
- Raina, R., Madhavan, A., and Ng, A. Y. (2009). Large-scale deep unsupervised learning using graphics processors. In L. Bottou and M. Littman, editors, *Proceedings of the Twenty-sixth International Conference on Machine Learning (ICML'09)*, pages 873–880, New York, NY, USA. ACM. [27](#), [449](#)
- Ramsey, F. P. (1926). Truth and probability. In R. B. Braithwaite, editor, *The Foundations of Mathematics and other Logical Essays*, chapter 7, pages 156–198. McMaster University Archive for the History of Economic Thought. [56](#)
- Ranzato, M. and Hinton, G. H. (2010). Modeling pixel means and covariances using factorized third-order Boltzmann machines. In *CVPR'2010*, pages 2551–2558. [682](#)
- Ranzato, M., Poultney, C., Chopra, S., and LeCun, Y. (2007a). Efficient learning of sparse representations with an energy-based model. In *NIPS'2006*. [14](#), [19](#), [510](#), [531](#), [533](#)
- Ranzato, M., Huang, F., Boureau, Y., and LeCun, Y. (2007b). Unsupervised learning of invariant feature hierarchies with applications to object recognition. In *Proceedings of the Computer Vision and Pattern Recognition Conference (CVPR'07)*. IEEE Press. [365](#)
- Ranzato, M., Boureau, Y., and LeCun, Y. (2008). Sparse feature learning for deep belief networks. In *NIPS'2007*. [510](#)
- Ranzato, M., Krizhevsky, A., and Hinton, G. E. (2010a). Factored 3-way restricted Boltzmann machines for modeling natural images. In *Proceedings of AISTATS 2010*. [680](#), [681](#)
- Ranzato, M., Mnih, V., and Hinton, G. (2010b). Generating more realistic images using gated MRFs. In *NIPS'2010*. [682](#), [683](#)
- Rao, C. (1945). Information and the accuracy attainable in the estimation of statistical parameters. *Bulletin of the Calcutta Mathematical Society*, **37**, 81–89. [135](#), [295](#)
- Rasmus, A., Valpola, H., Honkala, M., Berglund, M., and Raiko, T. (2015). Semi-supervised learning with ladder network. *arXiv preprint arXiv:1507.02672*. [429](#), [533](#)
- Recht, B., Re, C., Wright, S., and Niu, F. (2011). Hogwild: A lock-free approach to parallelizing stochastic gradient descent. In *NIPS'2011*. [450](#)
- Reichert, D. P., Seriès, P., and Storkey, A. J. (2011). Neuronal adaptation for sampling-based probabilistic inference in perceptual bistability. In *Advances in Neural Information Processing Systems*, pages 2357–2365. [668](#)

- Rezende, D. J., Mohamed, S., and Wierstra, D. (2014). Stochastic backpropagation and approximate inference in deep generative models. In *ICML'2014*. Preprint: arXiv:1401.4082. 655, 691, 698
- Rifai, S., Vincent, P., Muller, X., Glorot, X., and Bengio, Y. (2011a). Contractive auto-encoders: Explicit invariance during feature extraction. In *ICML'2011*. 524, 525, 526
- Rifai, S., Mesnil, G., Vincent, P., Muller, X., Bengio, Y., Dauphin, Y., and Glorot, X. (2011b). Higher order contractive auto-encoder. In *ECML PKDD*. 524, 525
- Rifai, S., Dauphin, Y., Vincent, P., Bengio, Y., and Muller, X. (2011c). The manifold tangent classifier. In *NIPS'2011*. 270, 271
- Rifai, S., Bengio, Y., Dauphin, Y., and Vincent, P. (2012). A generative process for sampling contractive auto-encoders. In *ICML'2012*. 713
- Ringach, D. and Shapley, R. (2004). Reverse correlation in neurophysiology. *Cognitive Science*, **28**(2), 147–166. 369
- Roberts, S. and Everson, R. (2001). *Independent component analysis: principles and practice*. Cambridge University Press. 496
- Robinson, A. J. and Fallside, F. (1991). A recurrent error propagation network speech recognition system. *Computer Speech and Language*, **5**(3), 259–274. 27, 462
- Rockafellar, R. T. (1997). Convex analysis. princeton landmarks in mathematics. 93
- Romero, A., Ballas, N., Ebrahimi Kahou, S., Chassang, A., Gatta, C., and Bengio, Y. (2015). Fitnets: Hints for thin deep nets. In *ICLR'2015*, arXiv:1412.6550. 326
- Rosen, J. B. (1960). The gradient projection method for nonlinear programming. part i. linear constraints. *Journal of the Society for Industrial and Applied Mathematics*, **8**(1), pp. 181–217. 93
- Rosenblatt, F. (1958). The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, **65**, 386–408. 14, 15, 27
- Rosenblatt, F. (1962). *Principles of Neurodynamics*. Spartan, New York. 15, 27
- Roweis, S. and Saul, L. K. (2000). Nonlinear dimensionality reduction by locally linear embedding. *Science*, **290**(5500). 163, 521
- Roweis, S., Saul, L., and Hinton, G. (2002). Global coordination of local linear models. In T. Dietterich, S. Becker, and Z. Ghahramani, editors, *Advances in Neural Information Processing Systems 14 (NIPS'01)*, Cambridge, MA. MIT Press. 492
- Rubin, D. B. *et al.* (1984). Bayesianly justifiable and relevant frequency calculations for the applied statistician. *The Annals of Statistics*, **12**(4), 1151–1172. 718

- Rumelhart, D., Hinton, G., and Williams, R. (1986a). Learning representations by back-propagating errors. *Nature*, **323**, 533–536. [14](#), [18](#), [23](#), [203](#), [225](#), [374](#), [479](#), [485](#)
- Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1986b). Learning internal representations by error propagation. In D. E. Rumelhart and J. L. McClelland, editors, *Parallel Distributed Processing*, volume 1, chapter 8, pages 318–362. MIT Press, Cambridge. [21](#), [27](#), [225](#)
- Rumelhart, D. E., McClelland, J. L., and the PDP Research Group (1986c). *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*. MIT Press, Cambridge. [17](#)
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A. C., and Fei-Fei, L. (2014a). ImageNet Large Scale Visual Recognition Challenge. [21](#)
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al. (2014b). Imagenet large scale visual recognition challenge. *arXiv preprint arXiv:1409.0575*. [28](#)
- Russel, S. J. and Norvig, P. (2003). *Artificial Intelligence: a Modern Approach*. Prentice Hall. [86](#)
- Rust, N., Schwartz, O., Movshon, J. A., and Simoncelli, E. (2005). Spatiotemporal elements of macaque V1 receptive fields. *Neuron*, **46**(6), 945–956. [368](#)
- Sainath, T., Mohamed, A., Kingsbury, B., and Ramabhadran, B. (2013). Deep convolutional neural networks for LVCSR. In *ICASSP 2013*. [463](#)
- Salakhutdinov, R. (2010). Learning in Markov random fields using tempered transitions. In Y. Bengio, D. Schuurmans, C. Williams, J. Lafferty, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22 (NIPS'09)*. [606](#)
- Salakhutdinov, R. and Hinton, G. (2009a). Deep Boltzmann machines. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, volume 5, pages 448–455. [24](#), [27](#), [532](#), [665](#), [668](#), [673](#), [674](#)
- Salakhutdinov, R. and Hinton, G. (2009b). Semantic hashing. In *International Journal of Approximate Reasoning*. [528](#)
- Salakhutdinov, R. and Hinton, G. E. (2007a). Learning a nonlinear embedding by preserving class neighbourhood structure. In *Proceedings of the Eleventh International Conference on Artificial Intelligence and Statistics (AISTATS'07)*, San Juan, Porto Rico. Omnipress. [530](#)
- Salakhutdinov, R. and Hinton, G. E. (2007b). Semantic hashing. In *SIGIR'2007*. [528](#)

- Salakhutdinov, R. and Hinton, G. E. (2008). Using deep belief nets to learn covariance kernels for Gaussian processes. In J. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems 20 (NIPS'07)*, pages 1249–1256, Cambridge, MA. MIT Press. [244](#)
- Salakhutdinov, R. and Larochelle, H. (2010). Efficient learning of deep Boltzmann machines. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics (AISTATS 2010), JMLR W&CP*, volume 9, pages 693–700. [655](#)
- Salakhutdinov, R. and Mnih, A. (2008). Probabilistic matrix factorization. In *NIPS'2008*. [482](#)
- Salakhutdinov, R. and Murray, I. (2008). On the quantitative analysis of deep belief networks. In W. W. Cohen, A. McCallum, and S. T. Roweis, editors, *Proceedings of the Twenty-fifth International Conference on Machine Learning (ICML'08)*, volume 25, pages 872–879. ACM. [631, 664](#)
- Salakhutdinov, R., Mnih, A., and Hinton, G. (2007). Restricted Boltzmann machines for collaborative filtering. In *ICML*. [482](#)
- Sanger, T. D. (1994). Neural network learning control of robot manipulators using gradually increasing task difficulty. *IEEE Transactions on Robotics and Automation*, **10**(3). [329](#)
- Saul, L. K. and Jordan, M. I. (1996). Exploiting tractable substructures in intractable networks. In D. Touretzky, M. Mozer, and M. Hasselmo, editors, *Advances in Neural Information Processing Systems 8 (NIPS'95)*. MIT Press, Cambridge, MA. [641](#)
- Saul, L. K., Jaakkola, T., and Jordan, M. I. (1996). Mean field theory for sigmoid belief networks. *Journal of Artificial Intelligence Research*, **4**, 61–76. [27, 695](#)
- Savich, A. W., Moussa, M., and Areibi, S. (2007). The impact of arithmetic representation on implementing mlp-bp on fpgas: A study. *Neural Networks, IEEE Transactions on*, **18**(1), 240–252. [454](#)
- Saxe, A. M., Koh, P. W., Chen, Z., Bhand, M., Suresh, B., and Ng, A. (2011). On random weights and unsupervised feature learning. In *Proc. ICML'2011*. ACM. [364](#)
- Saxe, A. M., McClelland, J. L., and Ganguli, S. (2013). Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. In *ICLR*. [285, 286, 303](#)
- Schaul, T., Antonoglou, I., and Silver, D. (2014). Unit tests for stochastic optimization. In *International Conference on Learning Representations*. [309](#)
- Schmidhuber, J. (1992). Learning complex, extended sequences using the principle of history compression. *Neural Computation*, **4**(2), 234–242. [401](#)
- Schmidhuber, J. (1996). Sequential neural text compression. *IEEE Transactions on Neural Networks*, **7**(1), 142–146. [480](#)

- Schmidhuber, J. (2012). Self-delimiting neural networks. *arXiv preprint arXiv:1210.0118*. 391
- Schölkopf, B. and Smola, A. J. (2002). *Learning with kernels: Support vector machines, regularization, optimization, and beyond*. MIT press. 705
- Schölkopf, B., Smola, A., and Müller, K.-R. (1998). Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, **10**, 1299–1319. 163, 521
- Schölkopf, B., Burges, C. J. C., and Smola, A. J. (1999). *Advances in Kernel Methods — Support Vector Learning*. MIT Press, Cambridge, MA. 18, 142
- Schölkopf, B., Janzing, D., Peters, J., Sgouritsa, E., Zhang, K., and Mooij, J. (2012). On causal and anticausal learning. In *ICML'2012*, pages 1255–1262. 548
- Schuster, M. (1999). On supervised learning from sequential data with applications for speech recognition. 189
- Schuster, M. and Paliwal, K. (1997). Bidirectional recurrent neural networks. *IEEE Transactions on Signal Processing*, **45**(11), 2673–2681. 396
- Schwenk, H. (2007). Continuous space language models. *Computer speech and language*, **21**, 492–518. 469
- Schwenk, H. (2010). Continuous space language models for statistical machine translation. *The Prague Bulletin of Mathematical Linguistics*, **93**, 137–146. 476
- Schwenk, H. (2014). Cleaned subset of WMT '14 dataset. 21
- Schwenk, H. and Bengio, Y. (1998). Training methods for adaptive boosting of neural networks. In M. Jordan, M. Kearns, and S. Solla, editors, *Advances in Neural Information Processing Systems 10 (NIPS'97)*, pages 647–653. MIT Press. 257
- Schwenk, H. and Gauvain, J.-L. (2002). Connectionist language modeling for large vocabulary continuous speech recognition. In *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 765–768, Orlando, Florida. 469
- Schwenk, H., Costa-jussà, M. R., and Fonollosa, J. A. R. (2006). Continuous space language models for the IWSLT 2006 task. In *International Workshop on Spoken Language Translation*, pages 166–173. 476
- Seide, F., Li, G., and Yu, D. (2011). Conversational speech transcription using context-dependent deep neural networks. In *Interspeech 2011*, pages 437–440. 23
- Sejnowski, T. (1987). Higher-order Boltzmann machines. In *AIP Conference Proceedings 151 on Neural Networks for Computing*, pages 398–403. American Institute of Physics Inc. 688

- Series, P., Reichert, D. P., and Storkey, A. J. (2010). Hallucinations in Charles Bonnet syndrome induced by homeostasis: a deep Boltzmann machine model. In *Advances in Neural Information Processing Systems*, pages 2020–2028. [668](#)
- Sermanet, P., Chintala, S., and LeCun, Y. (2012). Convolutional neural networks applied to house numbers digit classification. *CoRR*, [abs/1204.3968](#). [459](#)
- Sermanet, P., Kavukcuoglu, K., Chintala, S., and LeCun, Y. (2013). Pedestrian detection with unsupervised multi-stage feature learning. In *Proc. International Conference on Computer Vision and Pattern Recognition (CVPR'13)*. IEEE. [23](#), [200](#)
- Shilov, G. (1977). *Linear Algebra*. Dover Books on Mathematics Series. Dover Publications. [31](#)
- Siegelmann, H. (1995). Computation beyond the Turing limit. *Science*, **268**(5210), 545–548. [380](#)
- Siegelmann, H. and Sontag, E. (1991). Turing computability with neural nets. *Applied Mathematics Letters*, **4**(6), 77–80. [380](#)
- Siegelmann, H. T. and Sontag, E. D. (1995). On the computational power of neural nets. *Journal of Computer and Systems Sciences*, **50**(1), 132–150. [380](#), [406](#)
- Sietsma, J. and Dow, R. (1991). Creating artificial neural networks that generalize. *Neural Networks*, **4**(1), 67–79. [241](#)
- Simard, D., Steinkraus, P. Y., and Platt, J. C. (2003). Best practices for convolutional neural networks. In *ICDAR'2003*. [372](#)
- Simard, P. and Graf, H. P. (1994). Backpropagation without multiplication. In *Advances in Neural Information Processing Systems*, pages 232–239. [454](#)
- Simard, P., Victorri, B., LeCun, Y., and Denker, J. (1992). Tangent prop - A formalism for specifying selected invariances in an adaptive network. In *NIPS'1991*. [269](#), [271](#), [357](#)
- Simard, P. Y., LeCun, Y., and Denker, J. (1993). Efficient pattern recognition using a new transformation distance. In *NIPS'92*. [269](#)
- Simard, P. Y., LeCun, Y. A., Denker, J. S., and Victorri, B. (1998). Transformation invariance in pattern recognition — tangent distance and tangent propagation. *Lecture Notes in Computer Science*, **1524**. [269](#)
- Simons, D. J. and Levin, D. T. (1998). Failure to detect changes to people during a real-world interaction. *Psychonomic Bulletin & Review*, **5**(4), 644–649. [546](#)
- Simonyan, K. and Zisserman, A. (2015). Very deep convolutional networks for large-scale image recognition. In *ICLR*. [324](#)

- Sjöberg, J. and Ljung, L. (1995). Overtraining, regularization and searching for a minimum, with application to neural networks. *International Journal of Control*, **62**(6), 1391–1407. [249](#)
- Skinner, B. F. (1958). Reinforcement today. *American Psychologist*, **13**, 94–99. [329](#)
- Smolensky, P. (1986). Information processing in dynamical systems: Foundations of harmony theory. In D. E. Rumelhart and J. L. McClelland, editors, *Parallel Distributed Processing*, volume 1, chapter 6, pages 194–281. MIT Press, Cambridge. [574](#), [590](#), [658](#)
- Snoek, J., Larochelle, H., and Adams, R. P. (2012). Practical Bayesian optimization of machine learning algorithms. In *NIPS’2012*. [439](#)
- Socher, R., Huang, E. H., Pennington, J., Ng, A. Y., and Manning, C. D. (2011a). Dynamic pooling and unfolding recursive autoencoders for paraphrase detection. In *NIPS’2011*. [401](#), [403](#)
- Socher, R., Manning, C., and Ng, A. Y. (2011b). Parsing natural scenes and natural language with recursive neural networks. In *Proceedings of the Twenty-Eighth International Conference on Machine Learning (ICML’2011)*. [401](#)
- Socher, R., Pennington, J., Huang, E. H., Ng, A. Y., and Manning, C. D. (2011c). Semi-supervised recursive autoencoders for predicting sentiment distributions. In *EMNLP’2011*. [401](#)
- Socher, R., Perelygin, A., Wu, J. Y., Chuang, J., Manning, C. D., Ng, A. Y., and Potts, C. (2013a). Recursive deep models for semantic compositionality over a sentiment treebank. In *EMNLP’2013*. [401](#), [403](#)
- Socher, R., Ganjoo, M., Manning, C. D., and Ng, A. Y. (2013b). Zero-shot learning through cross-modal transfer. In *27th Annual Conference on Neural Information Processing Systems (NIPS 2013)*. [542](#)
- Sohl-Dickstein, J., Weiss, E. A., Maheswaranathan, N., and Ganguli, S. (2015). Deep unsupervised learning using nonequilibrium thermodynamics. [717](#), [718](#)
- Sohn, K., Zhou, G., and Lee, H. (2013). Learning and selecting features jointly with point-wise gated Boltzmann machines. In *ICML’2013*. [689](#)
- Solomonoff, R. J. (1989). A system for incremental learning based on algorithmic probability. [329](#)
- Sontag, E. D. (1998). VC dimension of neural networks. *NATO ASI Series F Computer and Systems Sciences*, **168**, 69–96. [550](#), [554](#)
- Sontag, E. D. and Sussman, H. J. (1989). Backpropagation can give rise to spurious local minima even for networks without hidden layers. *Complex Systems*, **3**, 91–106. [284](#)

- Sparkes, B. (1996). *The Red and the Black: Studies in Greek Pottery*. Routledge. 1
- Spitkovsky, V. I., Alshawi, H., and Jurafsky, D. (2010). From baby steps to leapfrog: how “less is more” in unsupervised dependency parsing. In *HLT’10*. 329
- Squire, W. and Trapp, G. (1998). Using complex variables to estimate derivatives of real functions. *SIAM Rev.*, **40**(1), 110–112. 442
- Srebro, N. and Shraibman, A. (2005). Rank, trace-norm and max-norm. In *Proceedings of the 18th Annual Conference on Learning Theory*, pages 545–560. Springer-Verlag. 239
- Srivastava, N. (2013). *Improving Neural Networks With Dropout*. Master’s thesis, U. Toronto. 538
- Srivastava, N. and Salakhutdinov, R. (2012). Multimodal learning with deep Boltzmann machines. In *NIPS’2012*. 544
- Srivastava, N., Salakhutdinov, R. R., and Hinton, G. E. (2013). Modeling documents with deep Boltzmann machines. *arXiv preprint arXiv:1309.6865*. 665
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2014). Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, **15**, 1929–1958. 257, 263, 264, 265, 674
- Srivastava, R. K., Greff, K., and Schmidhuber, J. (2015). Highway networks. *arXiv:1505.00387*. 327
- Steinkrau, D., Simard, P. Y., and Buck, I. (2005). Using GPUs for machine learning algorithms. *2013 12th International Conference on Document Analysis and Recognition*, **0**, 1115–1119. 448
- Stoyanov, V., Ropson, A., and Eisner, J. (2011). Empirical risk minimization of graphical model parameters given approximate inference, decoding, and model structure. In *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 15 of *JMLR Workshop and Conference Proceedings*, pages 725–733, Fort Lauderdale. Supplementary material (4 pages) also available. 676, 700
- Sukhbaatar, S., Szlam, A., Weston, J., and Fergus, R. (2015). Weakly supervised memory networks. *arXiv preprint arXiv:1503.08895*. 421
- Supancic, J. and Ramanan, D. (2013). Self-paced learning for long-term tracking. In *CVPR’2013*. 329
- Sussillo, D. (2014). Random walks: Training very deep nonlinear feed-forward networks with smart initialization. *CoRR*, **abs/1412.6558**. 290, 303, 305, 405
- Sutskever, I. (2012). *Training Recurrent Neural Networks*. Ph.D. thesis, Department of computer science, University of Toronto. 408, 415

- Sutskever, I. and Hinton, G. E. (2008). Deep narrow sigmoid belief networks are universal approximators. *Neural Computation*, **20**(11), 2629–2636. [695](#)
- Sutskever, I. and Tieleman, T. (2010). On the Convergence Properties of Contrastive Divergence. In Y. W. Teh and M. Titterington, editors, *Proc. of the International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 9, pages 789–795. [615](#)
- Sutskever, I., Hinton, G., and Taylor, G. (2009). The recurrent temporal restricted Boltzmann machine. In *NIPS'2008*. [688](#)
- Sutskever, I., Martens, J., and Hinton, G. E. (2011). Generating text with recurrent neural networks. In *ICML'2011*, pages 1017–1024. [480](#)
- Sutskever, I., Martens, J., Dahl, G., and Hinton, G. (2013). On the importance of initialization and momentum in deep learning. In *ICML*. [300](#), [408](#), [415](#)
- Sutskever, I., Vinyals, O., and Le, Q. V. (2014). Sequence to sequence learning with neural networks. In *NIPS'2014*, *arXiv:1409.3215*. [25](#), [101](#), [397](#), [411](#), [414](#), [477](#), [478](#)
- Sutton, R. and Barto, A. (1998). *Reinforcement Learning: An Introduction*. MIT Press. [106](#)
- Sutton, R. S., Mcallester, D., Singh, S., and Mansour, Y. (2000). Policy gradient methods for reinforcement learning with function approximation. In *NIPS'1999*, pages 1057–1063. MIT Press. [693](#)
- Swersky, K., Ranzato, M., Buchman, D., Marlin, B., and de Freitas, N. (2011). On autoencoders and score matching for energy based models. In *ICML'2011*. ACM. [516](#)
- Swersky, K., Snoek, J., and Adams, R. P. (2014). Freeze-thaw Bayesian optimization. *arXiv preprint arXiv:1406.3896*. [439](#)
- Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., and Rabinovich, A. (2014a). Going deeper with convolutions. Technical report, *arXiv:1409.4842*. [24](#), [27](#), [200](#), [257](#), [267](#), [327](#), [348](#)
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I. J., and Fergus, R. (2014b). Intriguing properties of neural networks. *ICLR*, **abs/1312.6199**. [267](#), [268](#), [270](#)
- Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., and Wojna, Z. (2015). Rethinking the Inception Architecture for Computer Vision. *ArXiv e-prints*. [244](#), [323](#)
- Taigman, Y., Yang, M., Ranzato, M., and Wolf, L. (2014). DeepFace: Closing the gap to human-level performance in face verification. In *CVPR'2014*. [100](#)
- Tandy, D. W. (1997). *Works and Days: A Translation and Commentary for the Social Sciences*. University of California Press. [1](#)

- Tang, Y. and Eliasmith, C. (2010). Deep networks for robust visual recognition. In *Proceedings of the 27th International Conference on Machine Learning, June 21-24, 2010, Haifa, Israel.* 241
- Tang, Y., Salakhutdinov, R., and Hinton, G. (2012). Deep mixtures of factor analysers. *arXiv preprint arXiv:1206.4635.* 492
- Taylor, G. and Hinton, G. (2009). Factored conditional restricted Boltzmann machines for modeling motion style. In L. Bottou and M. Littman, editors, *Proceedings of the Twenty-sixth International Conference on Machine Learning (ICML'09)*, pages 1025–1032, Montreal, Quebec, Canada. ACM. 688
- Taylor, G., Hinton, G. E., and Roweis, S. (2007). Modeling human motion using binary latent variables. In B. Schölkopf, J. Platt, and T. Hoffman, editors, *Advances in Neural Information Processing Systems 19 (NIPS'06)*, pages 1345–1352. MIT Press, Cambridge, MA. 687
- Teh, Y., Welling, M., Osindero, S., and Hinton, G. E. (2003). Energy-based models for sparse overcomplete representations. *Journal of Machine Learning Research*, 4, 1235–1260. 494
- Tenenbaum, J., de Silva, V., and Langford, J. C. (2000). A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500), 2319–2323. 163, 521, 536
- Theis, L., van den Oord, A., and Bethge, M. (2015). A note on the evaluation of generative models. *arXiv:1511.01844.* 699, 721
- Thompson, J., Jain, A., LeCun, Y., and Bregler, C. (2014). Joint training of a convolutional network and a graphical model for human pose estimation. In *NIPS'2014.* 361
- Thrun, S. (1995). Learning to play the game of chess. In *NIPS'1994.* 269
- Tibshirani, R. J. (1995). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society B*, 58, 267–288. 236
- Tieleman, T. (2008). Training restricted Boltzmann machines using approximations to the likelihood gradient. In W. W. Cohen, A. McCallum, and S. T. Roweis, editors, *Proceedings of the Twenty-fifth International Conference on Machine Learning (ICML'08)*, pages 1064–1071. ACM. 615
- Tieleman, T. and Hinton, G. (2009). Using fast weights to improve persistent contrastive divergence. In L. Bottou and M. Littman, editors, *Proceedings of the Twenty-sixth International Conference on Machine Learning (ICML'09)*, pages 1033–1040. ACM. 617
- Tipping, M. E. and Bishop, C. M. (1999). Probabilistic principal components analysis. *Journal of the Royal Statistical Society B*, 61(3), 611–622. 494

- Torralba, A., Fergus, R., and Weiss, Y. (2008). Small codes and large databases for recognition. In *Proceedings of the Computer Vision and Pattern Recognition Conference (CVPR'08)*, pages 1–8. [528](#)
- Touretzky, D. S. and Minton, G. E. (1985). Symbols among the neurons: Details of a connectionist inference architecture. In *Proceedings of the 9th International Joint Conference on Artificial Intelligence - Volume 1*, IJCAI'85, pages 238–243, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc. [17](#)
- Tu, K. and Honavar, V. (2011). On the utility of curricula in unsupervised learning of probabilistic grammars. In *IJCAI'2011*. [329](#)
- Turaga, S. C., Murray, J. F., Jain, V., Roth, F., Helmstaedter, M., Briggman, K., Denk, W., and Seung, H. S. (2010). Convolutional networks can learn to generate affinity graphs for image segmentation. *Neural Computation*, **22**(2), 511–538. [360](#)
- Turian, J., Ratinov, L., and Bengio, Y. (2010). Word representations: A simple and general method for semi-supervised learning. In *Proc. ACL'2010*, pages 384–394. [538](#)
- Töscher, A., Jahrer, M., and Bell, R. M. (2009). The BigChaos solution to the Netflix grand prize. [482](#)
- Uria, B., Murray, I., and Larochelle, H. (2013). Rnade: The real-valued neural autoregressive density-estimator. In *NIPS'2013*. [711](#)
- van den Oörd, A., Dieleman, S., and Schrauwen, B. (2013). Deep content-based music recommendation. In *NIPS'2013*. [483](#)
- van der Maaten, L. and Hinton, G. E. (2008). Visualizing data using t-SNE. *J. Machine Learning Res.*, **9**. [480, 522](#)
- Vanhoucke, V., Senior, A., and Mao, M. Z. (2011). Improving the speed of neural networks on CPUs. In *Proc. Deep Learning and Unsupervised Feature Learning NIPS Workshop*. [447, 455](#)
- Vapnik, V. N. (1982). *Estimation of Dependences Based on Empirical Data*. Springer-Verlag, Berlin. [114](#)
- Vapnik, V. N. (1995). *The Nature of Statistical Learning Theory*. Springer, New York. [114](#)
- Vapnik, V. N. and Chervonenkis, A. Y. (1971). On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and Its Applications*, **16**, 264–280. [114](#)
- Vincent, P. (2011). A connection between score matching and denoising autoencoders. *Neural Computation*, **23**(7). [516, 518, 714](#)

- Vincent, P. and Bengio, Y. (2003). Manifold Parzen windows. In *NIPS'2002*. MIT Press. **523**
- Vincent, P., Larochelle, H., Bengio, Y., and Manzagol, P.-A. (2008). Extracting and composing robust features with denoising autoencoders. In *ICML 2008*. **241, 518**
- Vincent, P., Larochelle, H., Lajoie, I., Bengio, Y., and Manzagol, P.-A. (2010). Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *J. Machine Learning Res.*, **11**. **518**
- Vincent, P., de Brébisson, A., and Bouthillier, X. (2015). Efficient exact gradient update for training deep networks with very large sparse targets. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 1108–1116. Curran Associates, Inc. **468**
- Vinyals, O., Kaiser, L., Koo, T., Petrov, S., Sutskever, I., and Hinton, G. (2014a). Grammar as a foreign language. Technical report, arXiv:1412.7449. **411**
- Vinyals, O., Toshev, A., Bengio, S., and Erhan, D. (2014b). Show and tell: a neural image caption generator. arXiv 1411.4555. **411**
- Vinyals, O., Fortunato, M., and Jaitly, N. (2015a). Pointer networks. *arXiv preprint arXiv:1506.03134*. **421**
- Vinyals, O., Toshev, A., Bengio, S., and Erhan, D. (2015b). Show and tell: a neural image caption generator. In *CVPR'2015*. arXiv:1411.4555. **102**
- Viola, P. and Jones, M. (2001). Robust real-time object detection. In *International Journal of Computer Vision*. **452**
- Visin, F., Kastner, K., Cho, K., Matteucci, M., Courville, A., and Bengio, Y. (2015). ReNet: A recurrent neural network based alternative to convolutional networks. *arXiv preprint arXiv:1505.00393*. **397**
- Von Melchner, L., Pallas, S. L., and Sur, M. (2000). Visual behaviour mediated by retinal projections directed to the auditory pathway. *Nature*, **404**(6780), 871–876. **16**
- Wager, S., Wang, S., and Liang, P. (2013). Dropout training as adaptive regularization. In *Advances in Neural Information Processing Systems 26*, pages 351–359. **264**
- Waibel, A., Hanazawa, T., Hinton, G. E., Shikano, K., and Lang, K. (1989). Phoneme recognition using time-delay neural networks. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, **37**, 328–339. **375, 456, 462**
- Wan, L., Zeiler, M., Zhang, S., LeCun, Y., and Fergus, R. (2013). Regularization of neural networks using dropconnect. In *ICML'2013*. **265**
- Wang, S. and Manning, C. (2013). Fast dropout training. In *ICML'2013*. **264**

- Wang, Z., Zhang, J., Feng, J., and Chen, Z. (2014a). Knowledge graph and text jointly embedding. In *Proc. EMNLP'2014*. 487
- Wang, Z., Zhang, J., Feng, J., and Chen, Z. (2014b). Knowledge graph embedding by translating on hyperplanes. In *Proc. AAAI'2014*. 487
- Warde-Farley, D., Goodfellow, I. J., Courville, A., and Bengio, Y. (2014). An empirical analysis of dropout in piecewise linear networks. In *ICLR'2014*. 261, 265, 266
- Wawrzynek, J., Asanovic, K., Kingsbury, B., Johnson, D., Beck, J., and Morgan, N. (1996). Spert-II: A vector microprocessor system. *Computer*, **29**(3), 79–86. 454
- Weaver, L. and Tao, N. (2001). The optimal reward baseline for gradient-based reinforcement learning. In *Proc. UAI'2001*, pages 538–545. 693
- Weinberger, K. Q. and Saul, L. K. (2004). Unsupervised learning of image manifolds by semidefinite programming. In *CVPR'2004*, pages 988–995. 163, 522
- Weiss, Y., Torralba, A., and Fergus, R. (2008). Spectral hashing. In *NIPS*, pages 1753–1760. 528
- Welling, M., Zemel, R. S., and Hinton, G. E. (2002). Self supervised boosting. In *Advances in Neural Information Processing Systems*, pages 665–672. 705
- Welling, M., Hinton, G. E., and Osindero, S. (2003a). Learning sparse topographic representations with products of Student-t distributions. In *NIPS'2002*. 682
- Welling, M., Zemel, R., and Hinton, G. E. (2003b). Self-supervised boosting. In S. Becker, S. Thrun, and K. Obermayer, editors, *Advances in Neural Information Processing Systems 15 (NIPS'02)*, pages 665–672. MIT Press. 626
- Welling, M., Rosen-Zvi, M., and Hinton, G. E. (2005). Exponential family harmoniums with an application to information retrieval. In L. Saul, Y. Weiss, and L. Bottou, editors, *Advances in Neural Information Processing Systems 17 (NIPS'04)*, volume 17, Cambridge, MA. MIT Press. 678
- Werbos, P. J. (1981). Applications of advances in nonlinear sensitivity analysis. In *Proceedings of the 10th IFIP Conference, 31.8 - 4.9, NYC*, pages 762–770. 225
- Weston, J., Bengio, S., and Usunier, N. (2010). Large scale image annotation: learning to rank with joint word-image embeddings. *Machine Learning*, **81**(1), 21–35. 403
- Weston, J., Chopra, S., and Bordes, A. (2014). Memory networks. *arXiv preprint arXiv:1410.3916*. 421, 488
- Widrow, B. and Hoff, M. E. (1960). Adaptive switching circuits. In *1960 IRE WESCON Convention Record*, volume 4, pages 96–104. IRE, New York. 15, 21, 24, 27

- Wikipedia (2015). List of animals by number of neurons — Wikipedia, the free encyclopedia. [Online; accessed 4-March-2015]. [24](#), [27](#)
- Williams, C. K. I. and Agakov, F. V. (2002). Products of Gaussians and Probabilistic Minor Component Analysis. *Neural Computation*, **14**(5), 1169–1182. [684](#)
- Williams, C. K. I. and Rasmussen, C. E. (1996). Gaussian processes for regression. In D. Touretzky, M. Mozer, and M. Hasselmo, editors, *Advances in Neural Information Processing Systems 8 (NIPS'95)*, pages 514–520. MIT Press, Cambridge, MA. [142](#)
- Williams, R. J. (1992). Simple statistical gradient-following algorithms connectionist reinforcement learning. *Machine Learning*, **8**, 229–256. [690](#), [691](#)
- Williams, R. J. and Zipser, D. (1989). A learning algorithm for continually running fully recurrent neural networks. *Neural Computation*, **1**, 270–280. [222](#)
- Wilson, D. R. and Martinez, T. R. (2003). The general inefficiency of batch training for gradient descent learning. *Neural Networks*, **16**(10), 1429–1451. [279](#)
- Wilson, J. R. (1984). Variance reduction techniques for digital simulation. *American Journal of Mathematical and Management Sciences*, **4**(3), 277—312. [692](#)
- Wiskott, L. and Sejnowski, T. J. (2002). Slow feature analysis: Unsupervised learning of invariances. *Neural Computation*, **14**(4), 715–770. [497](#)
- Wolpert, D. and MacReady, W. (1997). No free lunch theorems for optimization. *IEEE Transactions on Evolutionary Computation*, **1**, 67–82. [293](#)
- Wolpert, D. H. (1996). The lack of a priori distinction between learning algorithms. *Neural Computation*, **8**(7), 1341–1390. [116](#)
- Wu, R., Yan, S., Shan, Y., Dang, Q., and Sun, G. (2015). Deep image: Scaling up image recognition. arXiv:1501.02876. [450](#)
- Wu, Z. (1997). Global continuation for distance geometry problems. *SIAM Journal of Optimization*, **7**, 814–836. [328](#)
- Xiong, H. Y., Barash, Y., and Frey, B. J. (2011). Bayesian prediction of tissue-regulated splicing using RNA sequence and cellular context. *Bioinformatics*, **27**(18), 2554–2562. [264](#)
- Xu, K., Ba, J. L., Kiros, R., Cho, K., Courville, A., Salakhutdinov, R., Zemel, R. S., and Bengio, Y. (2015). Show, attend and tell: Neural image caption generation with visual attention. In *ICML'2015*, arXiv:1502.03044. [102](#), [411](#), [693](#)
- Yildiz, I. B., Jaeger, H., and Kiebel, S. J. (2012). Re-visiting the echo state property. *Neural networks*, **35**, 1–9. [407](#)

- Yosinski, J., Clune, J., Bengio, Y., and Lipson, H. (2014). How transferable are features in deep neural networks? In *NIPS'2014*. 324, 539
- Younes, L. (1998). On the convergence of Markovian stochastic algorithms with rapidly decreasing ergodicity rates. In *Stochastics and Stochastics Models*, pages 177–228. 615
- Yu, D., Wang, S., and Deng, L. (2010). Sequential labeling using deep-structured conditional random fields. *IEEE Journal of Selected Topics in Signal Processing*. 324
- Zaremba, W. and Sutskever, I. (2014). Learning to execute. arXiv 1410.4615. 330
- Zaremba, W. and Sutskever, I. (2015). Reinforcement learning neural Turing machines. *arXiv:1505.00521*. 422
- Zaslavsky, T. (1975). *Facing Up to Arrangements: Face-Count Formulas for Partitions of Space by Hyperplanes*. Number no. 154 in Memoirs of the American Mathematical Society. American Mathematical Society. 553
- Zeiler, M. D. and Fergus, R. (2014). Visualizing and understanding convolutional networks. In *ECCV'14*. 6
- Zeiler, M. D., Ranzato, M., Monga, R., Mao, M., Yang, K., Le, Q., Nguyen, P., Senior, A., Vanhoucke, V., Dean, J., and Hinton, G. E. (2013). On rectified linear units for speech processing. In *ICASSP 2013*. 462
- Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., and Torralba, A. (2015). Object detectors emerge in deep scene CNNs. *ICLR'2015*, arXiv:1412.6856. 554
- Zhou, J. and Troyanskaya, O. G. (2014). Deep supervised and convolutional generative stochastic network for protein secondary structure prediction. In *ICML'2014*. 717
- Zhou, Y. and Chellappa, R. (1988). Computation of optical flow using a neural network. In *Neural Networks, 1988., IEEE International Conference on*, pages 71–78. IEEE. 340
- Zöhrer, M. and Pernkopf, F. (2014). General stochastic networks for classification. In *NIPS'2014*. 717

# Index

- 0-1 loss, 104, 276
- Absolute value rectification, 192
- Accuracy, 426
- Activation function, 170
- Active constraint, 95
- AdaGrad, 307
- ADALINE, *see* adaptive linear element
- Adam, 308, 428
- Adaptive linear element, 15, 24, 27
- Adversarial example, 267
- Adversarial training, 268, 270, 533
- Affine, 110
- AIS, *see* annealed importance sampling
- Almost everywhere, 71
- Almost sure convergence, 130
- Ancestral sampling, 583, 598
- ANN, *see* Artificial neural network
- Annealed importance sampling, 628, 670, 719
- Approximate Bayesian computation, 718
- Approximate inference, 586
- Artificial intelligence, 1
- Artificial neural network, *see* Neural network
- ASR, *see* automatic speech recognition
- Asymptotically unbiased, 124
- Audio, 102, 361, 461
- Autoencoder, 4, 357, 505
- Automatic speech recognition, 461
- Back-propagation, 203
- Back-propagation through time, 385
- Backprop, *see* back-propagation
- Bag of words, 474
- Bagging, 255
- Batch normalization, 266, 428
- Bayes error, 117
- Bayes' rule, 70
- Bayesian hyperparameter optimization, 439
- Bayesian network, *see* directed graphical model
- Bayesian probability, 55
- Bayesian statistics, 135
- Belief network, *see* directed graphical model
- Bernoulli distribution, 62
- BFGS, 316
- Bias, 124, 229
- Bias parameter, 110
- Biased importance sampling, 596
- Bigram, 465
- Binary relation, 485
- Block Gibbs sampling, 602
- Boltzmann distribution, 573
- Boltzmann machine, 573, 656
- BPTT, *see* back-propagation through time
- Broadcasting, 34
- Burn-in, 600
- CAE, *see* contractive autoencoder
- Calculus of variations, 179
- Categorical distribution, *see* multinoulli distribution
- CD, *see* contrastive divergence
- Centering trick (DBM), 675
- Central limit theorem, 63
- Chain rule (calculus), 206
- Chain rule of probability, 59

- Chess, 2  
Chord, 582  
Chordal graph, 582  
Class-based language models, 466  
Classical dynamical system, 376  
Classification, 100  
Clique potential, *see* factor (graphical model)  
CNN, *see* convolutional neural network  
Collaborative Filtering, 481  
Collider, *see* explaining away  
Color images, 361  
Complex cell, 366  
Computational graph, 204  
Computer vision, 455  
Concept drift, 541  
Condition number, 279  
Conditional computation, *see* dynamic structure  
Conditional independence, xiii, 60  
Conditional probability, 59  
Conditional RBM, 687  
Connectionism, 17, 446  
Connectionist temporal classification, 463  
Consistency, 130, 516  
Constrained optimization, 93, 237  
Content-based addressing, 422  
Content-based recommender systems, 483  
Context-specific independence, 576  
Contextual bandits, 483  
Continuation methods, 328  
Contractive autoencoder, 524  
Contrast, 457  
Contrastive divergence, 291, 613, 674  
Convex optimization, 141  
Convolution, 331, 685  
Convolutional network, 16  
Convolutional neural network, 252, 331, 428, 463  
Coordinate descent, 322, 673  
Correlation, 61  
Cost function, *see* objective function  
Covariance, xiii, 61  
Covariance matrix, 62  
Coverage, 427  
Critical temperature, 606  
Cross-correlation, 333  
Cross-entropy, 75, 132  
Cross-validation, 122  
CTC, *see* connectionist temporal classification  
Curriculum learning, 329  
Curse of dimensionality, 154  
Cyc, 2  
D-separation, 575  
DAE, *see* denoising autoencoder  
Data generating distribution, 111, 131  
Data generating process, 111  
Data parallelism, 450  
Dataset, 105  
Dataset augmentation, 270, 460  
DBM, *see* deep Boltzmann machine  
DCGAN, 554, 555, 703  
Decision tree, 145, 551  
Decoder, 4  
Deep belief network, 27, 532, 634, 659, 662, 686, 694  
Deep Blue, 2  
Deep Boltzmann machine, 24, 27, 532, 634, 655, 659, 665, 674, 686  
Deep feedforward network, 167, 428  
Deep learning, 2, 5  
Denoising autoencoder, 513, 691  
Denoising score matching, 622  
Density estimation, 103  
Derivative, xiii, 83  
Design matrix, 106  
Detector layer, 340  
Determinant, xii  
Diagonal matrix, 41  
Differential entropy, 74, 649  
Dirac delta function, 65  
Directed graphical model, 77, 510, 566, 694  
Directional derivative, 85  
Discriminative fine-tuning, *see* supervised fine-tuning  
Discriminative RBM, 688  
Distributed representation, 17, 150, 549  
Domain adaptation, 539

- Dot product, 34, 141  
Double backprop, 270  
Doubly block circulant matrix, 334  
Dream sleep, 612, 655  
DropConnect, 265  
Dropout, 257, 428, 433, 434, 674, 691  
Dynamic structure, 451, 452
- E-step, 637  
Early stopping, 246, 249, 272, 273, 428  
EBM, *see* energy-based model  
Echo state network, 24, 27, 406  
Effective capacity, 114  
Eigendecomposition, 42  
Eigenvalue, 42  
Eigenvector, 42  
ELBO, *see* evidence lower bound  
Element-wise product, *see* Hadamard product, *see* Hadamard product  
EM, *see* expectation maximization  
Embedding, 519  
Empirical distribution, 66  
Empirical risk, 276  
Empirical risk minimization, 276  
Encoder, 4  
Energy function, 572  
Energy-based model, 572, 598, 656, 665  
Ensemble methods, 255  
Epoch, 247  
Equality constraint, 94  
Equivariance, 339  
Error function, *see* objective function  
ESN, *see* echo state network  
Euclidean norm, 39  
Euler-Lagrange equation, 649  
Evidence lower bound, 636, 663  
Example, 99  
Expectation, 60  
Expectation maximization, 637  
Expected value, *see* expectation  
Explaining away, 577, 634, 647  
Exploitation, 484  
Exploration, 484  
Exponential distribution, 65
- F-score, 426  
Factor (graphical model), 570  
Factor analysis, 493  
Factor graph, 582  
Factors of variation, 4  
Feature, 99  
Feature selection, 236  
Feedforward neural network, 167  
Fine-tuning, 324  
Finite differences, 442  
Forget gate, 306  
Forward propagation, 203  
Fourier transform, 361, 363  
Fovea, 367  
FPCD, 617  
Free energy, 574, 682  
Freebase, 486  
Frequentist probability, 55  
Frequentist statistics, 135  
Frobenius norm, 46  
Fully-visible Bayes network, 707  
Functional derivatives, 648  
FVBN, *see* fully-visible Bayes network
- Gabor function, 369  
GANs, *see* generative adversarial networks  
Gated recurrent unit, 428  
Gaussian distribution, *see* normal distribution  
Gaussian kernel, 142  
Gaussian mixture, 67, 188  
GCN, *see* global contrast normalization  
GeneOntology, 486  
Generalization, 110  
Generalized Lagrange function, *see* generalized Lagrangian  
Generalized Lagrangian, 94  
Generative adversarial networks, 691, 702  
Generative moment matching networks, 705  
Generator network, 695  
Gibbs distribution, 571  
Gibbs sampling, 584, 602  
Global contrast normalization, 457  
GPU, *see* graphics processing unit  
Gradient, 84

- Gradient clipping, 289, 417  
Gradient descent, 83, 85  
Graph, xii  
Graphical model, *see* structured probabilistic model  
Graphics processing unit, 447  
Greedy algorithm, 324  
Greedy layer-wise unsupervised pretraining, 531  
Greedy supervised pretraining, 324  
Grid search, 435  
Hadamard product, xii, 34  
Hard tanh, 196  
Harmonium, *see* restricted Boltzmann machine  
Harmony theory, 574  
Helmholtz free energy, *see* evidence lower bound  
Hessian, 223  
Hessian matrix, xiii, 87  
Heteroscedastic, 187  
Hidden layer, 6, 167  
Hill climbing, 86  
Hyperparameter optimization, 435  
Hyperparameters, 120, 433  
Hypothesis space, 112, 118  
i.i.d. assumptions, 111, 122, 267  
Identity matrix, 36  
ILSVRC, *see* ImageNet Large-Scale Visual Recognition Challenge  
ImageNet Large-Scale Visual Recognition Challenge, 23  
Immorality, 580  
Importance sampling, 595, 627, 700  
Importance weighted autoencoder, 700  
Independence, xiii, 60  
Independent and identically distributed, *see* i.i.d. assumptions  
Independent component analysis, 494  
Independent subspace analysis, 496  
Inequality constraint, 94  
Inference, 565, 586, 634, 636, 638, 641, 651, 653  
Information retrieval, 528  
Initialization, 301  
Integral, xiii  
Invariance, 343  
Isotropic, 65  
Jacobian matrix, xiii, 72, 86  
Joint probability, 57  
*k*-means, 365, 549  
*k*-nearest neighbors, 143, 551  
Karush-Kuhn-Tucker conditions, 95, 237  
Karush–Kuhn–Tucker, 94  
Kernel (convolution), 332, 333  
Kernel machine, 551  
Kernel trick, 141  
KKT, *see* Karush–Kuhn–Tucker  
KKT conditions, *see* Karush–Kuhn–Tucker conditions  
KL divergence, *see* Kullback–Leibler divergence  
Knowledge base, 2, 486  
Krylov methods, 224  
Kullback–Leibler divergence, xiii, 74  
Label smoothing, 243  
Lagrange multipliers, 94, 649  
Lagrangian, *see* generalized Lagrangian  
LAPGAN, 704  
Laplace distribution, 65, 499, 500  
Latent variable, 67  
Layer (neural network), 167  
LCN, *see* local contrast normalization  
Leaky ReLU, 192  
Leaky units, 409  
Learning rate, 85  
Line search, 85, 86, 93  
Linear combination, 37  
Linear dependence, 38  
Linear factor models, 492  
Linear regression, 107, 110, 140  
Link prediction, 487  
Lipschitz constant, 92  
Lipschitz continuous, 92  
Liquid state machine, 406

- Local conditional probability distribution, 567  
Local contrast normalization, 459  
Logistic regression, 3, 140, 140  
Logistic sigmoid, 7, 67  
Long short-term memory, 18, 25, 306, 411, 428  
Loop, 582  
Loopy belief propagation, 588  
Loss function, *see* objective function  
 $L^p$  norm, 39  
LSTM, *see* long short-term memory  
M-step, 637  
Machine learning, 2  
Machine translation, 101  
Main diagonal, 33  
Manifold, 160  
Manifold hypothesis, 161  
Manifold learning, 161  
Manifold tangent classifier, 270  
MAP approximation, 138, 508  
Marginal probability, 58  
Markov chain, 598  
Markov chain Monte Carlo, 598  
Markov network, *see* undirected model  
Markov random field, *see* undirected model  
Matrix, xi, xii, 32  
Matrix inverse, 36  
Matrix product, 34  
Max norm, 40  
Max pooling, 340  
Maximum likelihood, 131  
Maxout, 192, 428  
MCMC, *see* Markov chain Monte Carlo  
Mean field, 641, 642, 674  
Mean squared error, 108  
Measure theory, 71  
Measure zero, 71  
Memory network, 419, 421  
Method of steepest descent, *see* gradient descent  
Minibatch, 279  
Missing inputs, 100  
Mixing (Markov chain), 604  
Mixture density networks, 188  
Mixture distribution, 66  
Mixture model, 188, 513  
Mixture of experts, 453, 551  
MLP, *see* multilayer perception  
MNIST, 21, 22, 674  
Model averaging, 255  
Model compression, 451  
Model identifiability, 284  
Model parallelism, 450  
Moment matching, 705  
Moore-Penrose pseudoinverse, 45, 240  
Moralized graph, 580  
MP-DBM, *see* multi-prediction DBM  
MRF (Markov Random Field), *see* undirected model  
MSE, *see* mean squared error  
Multi-modal learning, 542  
Multi-prediction DBM, 676  
Multi-task learning, 245, 541  
Multilayer perception, 5  
Multilayer perceptron, 27  
Multinomial distribution, 62  
Multinoulli distribution, 62  
 $n$ -gram, 464  
NADE, 710  
Naive Bayes, 3  
Nat, 73  
Natural image, 562  
Natural language processing, 464  
Nearest neighbor regression, 115  
Negative definite, 89  
Negative phase, 473, 609, 611  
Neocognitron, 16, 24, 27, 368  
Nesterov momentum, 300  
Netflix Grand Prize, 256, 482  
Neural language model, 466, 479  
Neural network, 13  
Neural Turing machine, 421  
Neuroscience, 15  
Newton's method, 89, 310  
NLM, *see* neural language model  
NLP, *see* natural language processing  
No free lunch theorem, 116

- Noise-contrastive estimation, 623  
Non-parametric model, 114  
Norm, xiv, 39  
Normal distribution, 63, 64, 125  
Normal equations, 109, 109, 112, 234  
Normalized initialization, 303  
Numerical differentiation, *see* finite differences  
  
Object detection, 456  
Object recognition, 456  
Objective function, 82  
OMP- $k$ , *see* orthogonal matching pursuit  
One-shot learning, 541  
Operation, 204  
Optimization, 80, 82  
Orthodox statistics, *see* frequentist statistics  
Orthogonal matching pursuit, 27, 254  
Orthogonal matrix, 42  
Orthogonality, 41  
Output layer, 167  
  
Parallel distributed processing, 17  
Parameter initialization, 301, 408  
Parameter sharing, 251, 336, 374, 376, 389  
Parameter tying, *see* Parameter sharing  
Parametric model, 114  
Parametric ReLU, 192  
Partial derivative, 84  
Partition function, 571, 608, 671  
PCA, *see* principal components analysis  
PCD, *see* stochastic maximum likelihood  
Perceptron, 15, 27  
Persistent contrastive divergence, *see* stochastic maximum likelihood  
Perturbation analysis, *see* reparametrization trick  
Point estimator, 122  
Policy, 483  
Pooling, 331, 685  
Positive definite, 89  
Positive phase, 473, 609, 611, 658, 670  
Precision, 426  
Precision (of a normal distribution), 63, 65  
Predictive sparse decomposition, 526  
  
Preprocessing, 456  
Pretraining, 324, 531  
Primary visual cortex, 366  
Principal components analysis, 48, 146–148, 493, 634  
Prior probability distribution, 135  
Probabilistic max pooling, 685  
Probabilistic PCA, 493, 494, 635  
Probability density function, 58  
Probability distribution, 56  
Probability mass function, 56  
Probability mass function estimation, 103  
Product of experts, 573  
Product rule of probability, *see* chain rule of probability  
PSD, *see* predictive sparse decomposition  
Pseudolikelihood, 618  
  
Quadrature pair, 370  
Quasi-Newton condition, 316  
Quasi-Newton methods, 316  
  
Radial basis function, 196  
Random search, 437  
Random variable, 56  
Ratio matching, 621  
RBF, 196  
RBM, *see* restricted Boltzmann machine  
Recall, 426  
Receptive field, 338  
Recommender Systems, 481  
Rectified linear unit, 171, 192, 428, 510  
Recurrent network, 27  
Recurrent neural network, 379  
Regression, 101  
Regularization, 120, 120, 177, 228, 433  
Regularizer, 119  
REINFORCE, 691  
Reinforcement learning, 25, 106, 483, 691  
Relational database, 486  
Reparametrization trick, 690  
Representation learning, 3  
Representational capacity, 114  
Restricted Boltzmann machine, 357, 462, 482, 590, 634, 658, 659, 674, 678,

- 680, 683, 685
- Ridge regression, *see* weight decay
- Risk, 275
- RNN-RBM, 688
- Saddle points, 285
- Sample mean, 125
- Scalar, xi, xii, 31
- Score matching, 516, 620
- Secant condition, 316
- Second derivative, 86
- Second derivative test, 89
- Self-information, 73
- Semantic hashing, 528
- Semi-supervised learning, 244
- Separable convolution, 363
- Separation (probabilistic modeling), 575
- Set, xii
- SGD, *see* stochastic gradient descent
- Shannon entropy, xiii, 73
- Shortlist, 469
- Sigmoid, xiv, *see* logistic sigmoid
- Sigmoid belief network, 27
- Simple cell, 366
- Singular value, *see* singular value decomposition
- Singular value decomposition, 44, 148, 482
- Singular vector, *see* singular value decomposition
- Slow feature analysis, 496
- SML, *see* stochastic maximum likelihood
- Softmax, 183, 421, 453
- Softplus, xiv, 68, 196
- Spam detection, 3
- Sparse coding, 322, 357, 499, 634, 694
- Sparse initialization, 304, 408
- Sparse representation, 146, 226, 253, 508, 559
- Spearmint, 439
- Spectral radius, 407
- Speech recognition, *see* automatic speech recognition
- Sphering, *see* whitening
- Spike and slab restricted Boltzmann machine, 683
- SPN, *see* sum-product network
- Square matrix, 38
- ssRBM, *see* spike and slab restricted Boltzmann machine
- Standard deviation, 61
- Standard error, 127
- Standard error of the mean, 128, 278
- Statistic, 122
- Statistical learning theory, 110
- Steepest descent, *see* gradient descent
- Stochastic back-propagation, *see* reparametrization trick
- Stochastic gradient descent, 15, 150, 279, 294, 674
- Stochastic maximum likelihood, 615, 674
- Stochastic pooling, 265
- Structure learning, 585
- Structured output, 101, 687
- Structured probabilistic model, 77, 561
- Sum rule of probability, 58
- Sum-product network, 556
- Supervised fine-tuning, 532, 664
- Supervised learning, 105
- Support vector machine, 140
- Surrogate loss function, 276
- SVD, *see* singular value decomposition
- Symmetric matrix, 41, 43
- Tangent distance, 269
- Tangent plane, 519
- Tangent prop, 269
- TDNN, *see* time-delay neural network
- Teacher forcing, 383, 384
- Tempering, 606
- Template matching, 141
- Tensor, xi, xii, 33
- Test set, 110
- Tikhonov regularization, *see* weight decay
- Tiled convolution, 353
- Time-delay neural network, 369, 375
- Toeplitz matrix, 334
- Topographic ICA, 496
- Trace operator, 46
- Training error, 110
- Transcription, 101

- Transfer learning, 539  
Transpose, xii, 33  
Triangle inequality, 39  
Triangulated graph, *see* chordal graph  
Trigram, 465
- Unbiased, 124  
Undirected graphical model, 77, 510  
Undirected model, 569  
Uniform distribution, 57  
Unigram, 465  
Unit norm, 41  
Unit vector, 41  
Universal approximation theorem, 197  
Universal approximator, 556  
Unnormalized probability distribution, 570  
Unsupervised learning, 105, 146  
Unsupervised pretraining, 462, 531
- V-structure, *see* explaining away  
V1, 366  
VAE, *see* variational autoencoder  
Vapnik-Chervonenkis dimension, 114  
Variance, xiii, 61, 229  
Variational autoencoder, 691, 698  
Variational derivatives, *see* functional derivatives  
Variational free energy, *see* evidence lower bound  
VC dimension, *see* Vapnik-Chervonenkis dimension  
Vector, xi, xii, 32  
Virtual adversarial examples, 268  
Visible layer, 6  
Volumetric data, 361
- Wake-sleep, 654, 663  
Weight decay, 118, 177, 231, 434  
Weight space symmetry, 284  
Weights, 15, 107  
Whitening, 458  
Wikibase, 486  
**Wikibase**, 486  
Word embedding, 467  
Word-sense disambiguation, 487