

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} . 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 section 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