

Greedy Layer-wise Unsupervised Pretraining

Sargur N. Srihari
srihari@cedar.buffalo.edu

Topics in Representation Learning

- Overview
 1. Greedy Layer-Wise Unsupervised Pretraining
 2. Transfer Learning and Domain Adaptation
 3. Semi-supervised Disentangling of Causal Factors
 4. Distributed Representation
 5. Exponential Gains from depth
 6. Providing Clues to Discover Underlying Causes

Pre-training and fine tuning

- Using dataset A train model M
- Pre-training:
 - You have a dataset B
 - Before training the model, initialize some of the parameters of M with model trained on A
- Fine-tuning:
 - You train M on B
- This is one form of transfer learning

Unsupervised Pretraining

- Unsupervised learning played key historical role in revival of deep neural networks
 - Enabling training a deep supervised network without requiring architectural specializations such as convolution or recurrence
- We call this procedure *unsupervised pretraining*
- Or more precisely *greedy layer-wise unsupervised pretraining*

Greedy Algorithm

- Greedy algorithms break a problem into many components, then solve for the optimal version of each component in isolation
- Unfortunately, combining the individually optimal components is not guaranteed to yield an optimal complete solution

Greedy Layer-wise Unsupervised Pretraining

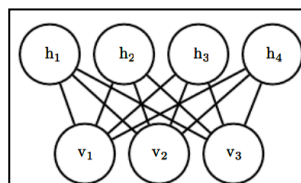
- A representation learned for one task
 - unsupervised learning, that captures the shape of the input distribution
- Is used for another task
 - supervised learning with the same input domain
- Greedy layer-wise pre-training relies on a single-layer representation learning

Single-layer representation learning

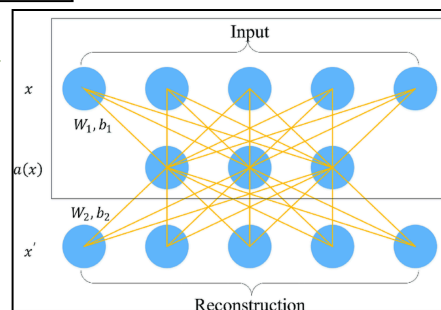
- We need a single-layer representation learning algorithm, such as:

- An *RBM*

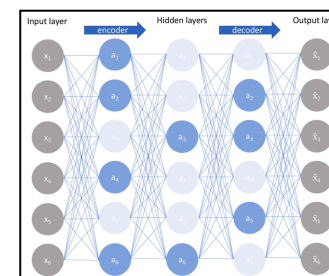
- (a Markov network)



- A single-layer *autoencoder*



- A sparse coding model



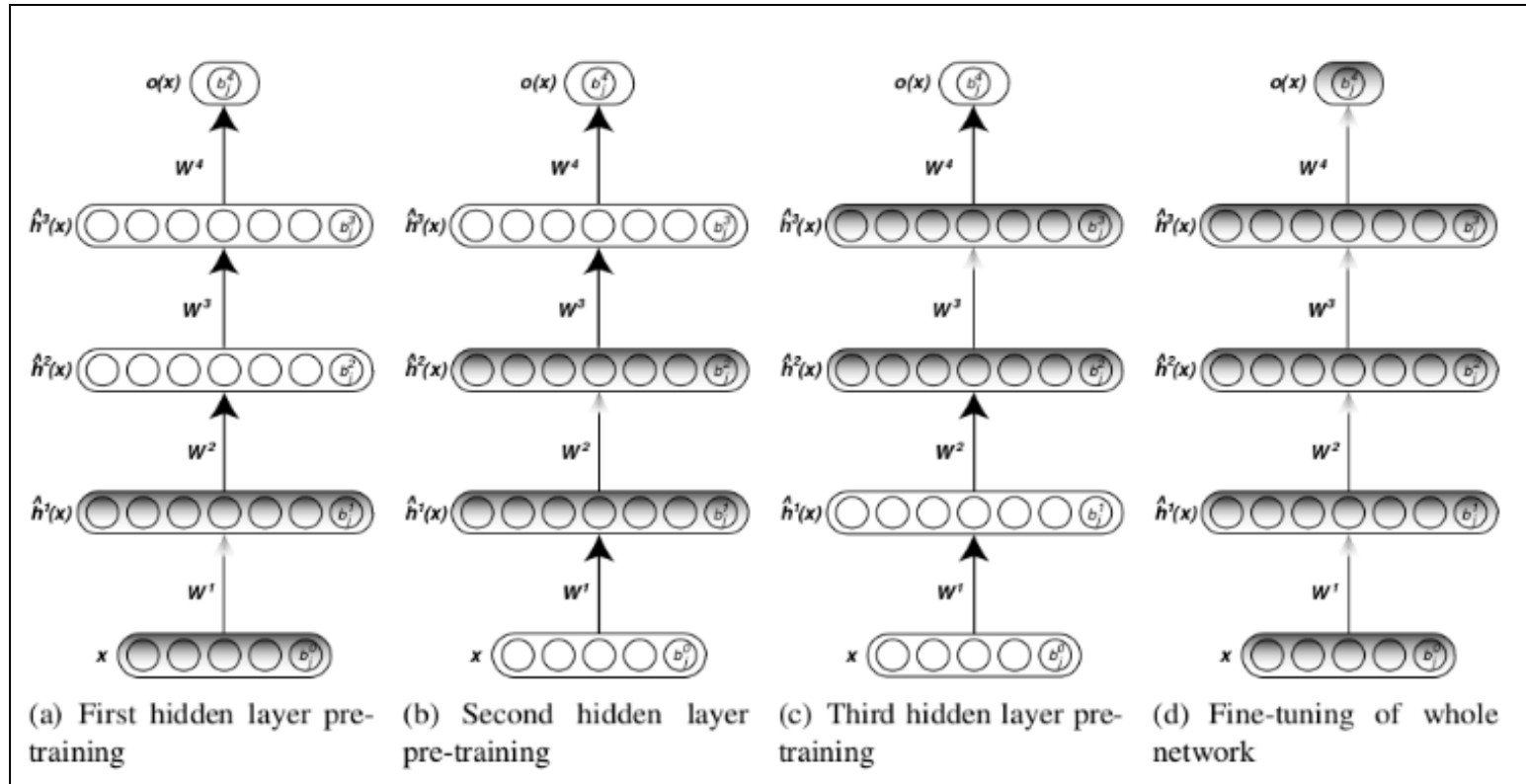
- Or another model that learns latent representations

Single layer Pretraining

- Each layer pretrained using unsupervised learning
 - Taking the output of the previous layer and producing as output a new representation of data,
 - Whose distribution (or relation to categories) is simpler

Training a 4-layer network

- Pairs of layers active in each stage



Formal Algorithm

Algorithm: Greedy Layer-wise Unsupervised Pretraining Protocol

- Given unsupervised feature learning algorithm \mathcal{L}
 - Which takes as input a training set of examples and returns an encoder or feature function f
- Raw input data is X , with one row per example, $f^{(1)}(X)$ is output of the first stage encoder on X
- In the case where fine tuning is performed we use
 - a learner \mathcal{T} which takes an initial function f , input examples X (and in the supervised fine-tuning case, associated targets Y) and returns a tuned function.
 - The no of stages is m

```

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

```

History: layer-wise unsupervised

- Unsupervised greedy layer-wise training
 - was used to sidestep difficulty of training layers of a deep neural net for a supervised task
 - Origins in Neocognitron (Fukushima, 1975)
 - Deep learning renaissance of 2006 began with
 - Greedy learning to find initialization for all layers
 - Useful for fully connected architectures
 - Earlier, only deep CNNs or depth resulting from recurrence were feasible to train
- Today greedy layer-wise pretraining is not required to train fully connected deep networks

Greedy pretraining terminology

- Greedy layer-wise pretraining
 - Greedy because
 - It is a greedy algorithm that optimizes each piece of the solution independently
 - One piece at a time rather than jointly
 - Layer-wise because
 - Independent pieces are the layers of the network
 - Training proceeds one layer at a time
 - Training the k^{th} layer while previous ones are fixed
 - Pretraining because
 - It is only a first step before applying a joint training algorithm is applied to *fine-tune* all layers together

When/why does pretraining work?

- Greedy layer-wise unsupervised pretraining can yield substantial improvements for classification
 - However it is sometimes *harmful*
- So it is important to learn as to when it works to determine whether it is applicable to a particular task
- This discussion pertains only to greedy unsupervised pretraining. There are other semi-supervised learning paradigms

Unsupervised pretraining combines two ideas

1. Initial parameters have a regularizing effect
 - i.e., approach one local minimum over another
 - But local minima no longer considered serious
2. Learning about input distribution can help to learn about the mapping from inputs to outputs
 - Learns that cars and motorcycles have wheels
 - The representation for wheels is useful for the supervised learner

Learning trajectories

- Each point refers to a neural network at a particular time in its training
- Pretraining accesses new part of space:
 - With pretraining: halt in one region of function space
 - Without pretraining: another region

Visualization of functions projected into 2d space.
(Each function is an infinite-dimensional vector that associates every input x with output y).
Color indicates time.
Area where pretrained networks arrive is smaller

