

Greedy Layer-Wise Training of Deep Networks

**Yoshua Bengio, Pascal Lamblin, Dan Popovici,
Hugo Larochelle**

December 5th 2006

**Thanks to: Yann Le Cun, Geoffrey Hinton,
Olivier Delalleau, Nicolas Le Roux**

- Motivation: AI \Rightarrow highly-varying functions
 \Rightarrow non-local + **deep architecture**
- Principle of **greedy layer-wise unsupervised initialization**
- Deep Belief Networks
- Deep Multi-layer Neural Networks
- Experimental study: **why this principle works**
- Extensions of Deep Belief Networks

Grand Goal: AI

- **Goal: using ML to reach AI**
- *AI tasks*: visual and auditory perception, language understanding, intelligent control, long-term prediction, understanding of high-level abstractions...
- **Remains elusive!** (and mostly forgotten?)
- 3 considerations:
 - computational efficiency
 - statistical efficiency
 - human-labor efficiency
- Here: focus on algorithms associated with broad priors (i.e., non-parametric) with the hope of discovering principles applicable to vast array of tasks within AI, with no need of hand-crafted solutions for each particular task.

Grand Goal: AI

- Goal: using ML to reach AI
- AI tasks: visual and auditory perception, language understanding, intelligent control, long-term prediction, understanding of high-level abstractions...
- Remains elusive! (and mostly forgotten?)
- 3 considerations:
 - computational efficiency
 - statistical efficiency
 - human-labor efficiency
- Here: focus on algorithms associated with broad priors (i.e., non-parametric) with the hope of discovering principles applicable to vast array of tasks within AI, with no need of hand-crafted solutions for each particular task.

Depth of Architectures

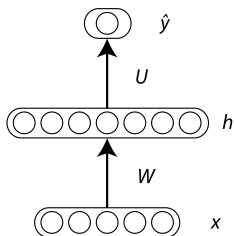
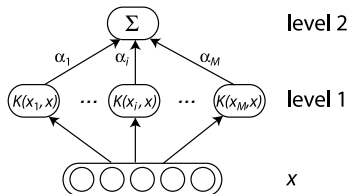
Depth = number of **levels** of composition of adaptable elements:

- kernel machines: **shallow**
- boosting: generally **shallow**
- multi-layer neural networks: **usually shallow, can be deep?**
- decision trees: deep but local estimators (curse of dim.)
- parametric graphical models: human-labor intensive

Non-parametric ones can theoretically approximate any continuous function.

But **how efficiently?**

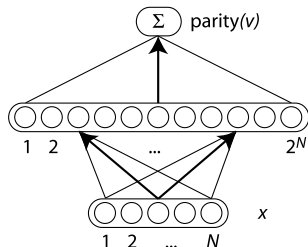
(computational, statistical)



Inefficiency of Shallow Architectures

Shallow architectures may need huge number of elements.

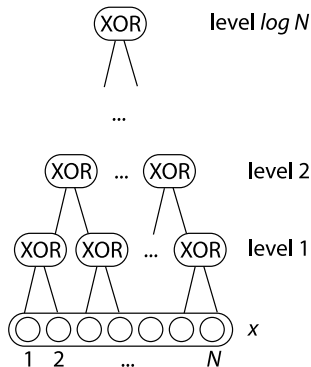
Worst-case can be **exponentially bad**. Many theorems and examples from boolean circuits theory (Hastad 1987): multiplier circuits, parity, etc..



Very fat shallow circuit

\Rightarrow many adjustable elements

\Rightarrow many examples needed



Number of levels should increase with complexity, nb inputs, maybe logarithmic.

Curse of Dimensionality

(Bengio et al 2006):

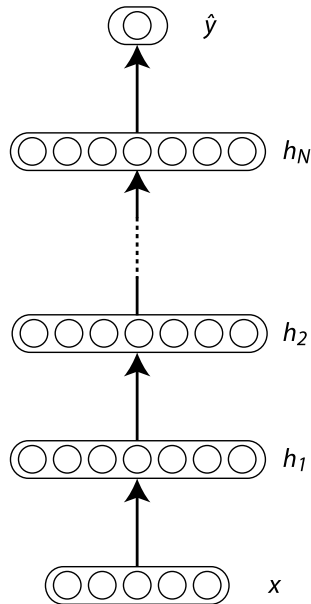
Local kernel machines (= **pattern matchers**) and decision trees **partition the space** and may need

exponential nb of units, i.e. of examples

inefficient at representing **highly-varying functions**, which may otherwise have a compact representation.

Optimization of Deep Architectures

- What **deep** architectures are known? various kinds of multi-layer neural networks with many layers.
- Except for a very special kind of architectures for machine vision (convolutional networks), deep architectures have been neglected in machine learning.
- Why? **Training gets stuck in mediocre solutions** (Tesauro 92).
- Credit assignment problem?
- No hope?



Greedy Learning of Multiple Levels of Abstractions

- **Greedily learning simple things first, higher-level abstractions on top of lower-level ones seems like a good strategy.**

Greedy Learning of Multiple Levels of Abstractions

- **Greedly learning simple things first, higher-level abstractions on top of lower-level ones seems like a good strategy.**
- **Implicit prior:** restrict to functions that
 - 1 can be represented as a composition of simpler ones such that
 - 2 the simpler ones can be learned first (i.e., are also good models of the data).

Greedy Learning of Multiple Levels of Abstractions

- **Greedly learning simple things first, higher-level abstractions on top of lower-level ones seems like a good strategy.**
- **Implicit prior:** restrict to functions that
 - ① can be represented as a composition of simpler ones such that
 - ② the simpler ones can be learned first (i.e., are also good models of the data).
- Coherent with psychological literature (Piaget 1952).
We learn baby math before arithmetic before algebra before differential equations . . .
- Also some evidence from neurobiology: (Guillery 2005) *“Is postnatal neocortical maturation hierarchical?”*.

Deep Belief Networks

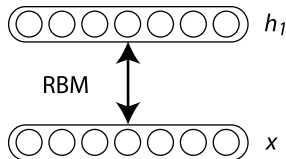
Hinton et al (2006) just introduced a deep network model that provides more evidence that this direction is worthwhile:

- beats state-of-the-art statistical learning in experiments on a large machine learning benchmark task (MNIST)

Deep Belief Networks

Hinton et al (2006) just introduced a deep network model that provides more evidence that this direction is worthwhile:

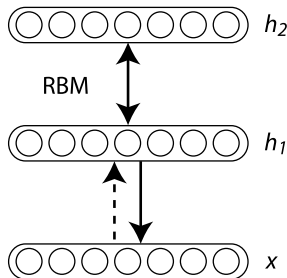
- beats state-of-the-art statistical learning in experiments on a large machine learning benchmark task (MNIST)
- Each layer tries to model distribution of its input (unsupervised training as Restricted Boltzmann Machine)
- H = hidden causes,
 $P(h|v)$ = representation of v .



Deep Belief Networks

Hinton et al (2006) just introduced a deep network model that provides more evidence that this direction is worthwhile:

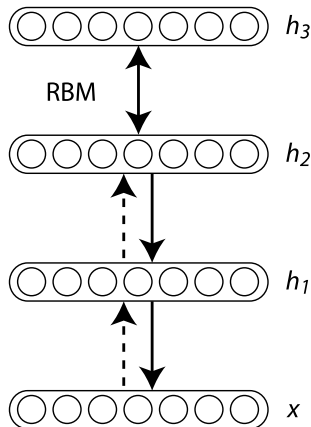
- **beats state-of-the-art statistical learning in experiments on a large machine learning benchmark task (MNIST)**
- Each layer tries to model distribution of its input (unsupervised training as Restricted Boltzmann Machine)
- H = hidden causes,
 $P(h|v)$ = representation of v .
- **Unsupervised greedy layer-wise training = initialization**, replaces traditional random initialization of multi-layer networks



Deep Belief Networks

Hinton et al (2006) just introduced a deep network model that provides more evidence that this direction is worthwhile:

- beats state-of-the-art statistical learning in experiments on a large machine learning benchmark task (MNIST)
- Each layer tries to model distribution of its input (unsupervised training as Restricted Boltzmann Machine)
- H = hidden causes,
 $P(h|v)$ = representation of v .
- Unsupervised greedy layer-wise training = **initialization**, replaces traditional random initialization of multi-layer networks

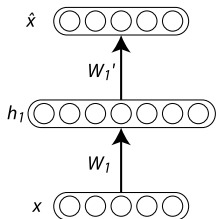


Greedy Layer-Wise Initialization

- The principle of greedy layer-wise initialization proposed by Hinton can be generalized to other algorithms.
- Initialize each layer of a deep multi-layer feedforward neural net as an autoassociator for the output of previous layer.
- Find W which minimizes cross-entropy loss in predicting x from $\hat{x} = \text{sigm}(W' \text{sigm}(Wx))$.
- Feed its hidden activations as input to next layer.
- Small weights (weight decay or stochastic gradient) prevent autoassociator from learning the identity.

Greedy Layer-Wise Initialization

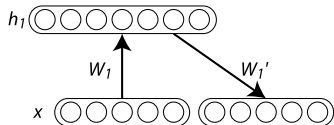
- The principle of greedy layer-wise initialization proposed by Hinton can be generalized to other algorithms.
- Initialize each layer of a **deep multi-layer feedforward neural net** as an **autoassociator** for the output of previous layer.
- Find W which minimizes cross-entropy loss in predicting x from $\hat{x} = \text{sigm}(W' \text{sigm}(Wx))$.



- Feed its hidden activations as input to next layer.
- Small weights (weight decay or stochastic gradient) prevent autoassociator from learning the identity.

Greedy Layer-Wise Initialization

- The principle of greedy layer-wise initialization proposed by Hinton can be generalized to other algorithms.
- Initialize each layer of a deep multi-layer feedforward neural net as an autoassociator for the output of previous layer.
- Find W which minimizes cross-entropy loss in predicting x from $\hat{x} = \text{sigm}(W' \text{sigm}(Wx))$.



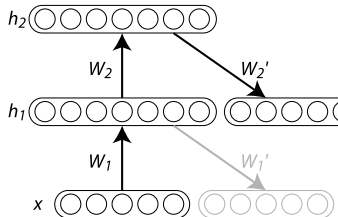
- Feed its hidden activations as input to next layer.
- Small weights (weight decay or stochastic gradient) prevent autoassociator from learning the identity.

Greedy Layer-Wise Initialization

- The principle of greedy layer-wise initialization proposed by Hinton can be generalized to other algorithms.

- Initialize each layer of a **deep multi-layer feedforward neural net** as an **autoassociator** for the output of previous layer.

- Find W which minimizes cross-entropy loss in predicting x from $\hat{x} = \text{sigm}(W' \text{sigm}(Wx))$.

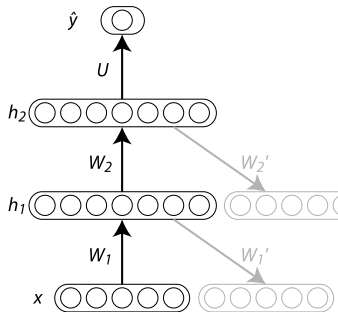


- Feed its hidden activations as input to next layer.
- Small weights (weight decay or stochastic gradient) prevent autoassociator from learning the identity.

Greedy Layer-Wise Initialization

- The principle of greedy layer-wise initialization proposed by Hinton can be generalized to other algorithms.

- Initialize each layer of a **deep multi-layer feedforward neural net** as an **autoassociator** for the output of previous layer.
- Find W which minimizes cross-entropy loss in predicting x from $\hat{x} = \text{sigm}(W' \text{sigm}(Wx))$.

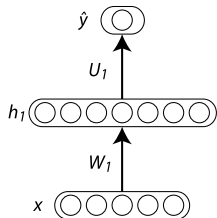


- Feed its hidden activations as input to next layer.
- Small weights (weight decay or stochastic gradient) prevent autoassociator from learning the identity.

Greedy Supervised Layer-Wise Initialization

Why not use supervised learning at each stage?

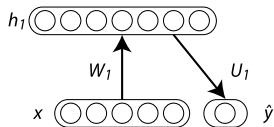
- Each layer is trained as the hidden layer of a supervised 2-layer neural net.
- After training the 2-layer neural net, discard output layer;
- Propagate data through new hidden layer, train another layer, etc.



Greedy Supervised Layer-Wise Initialization

Why not use supervised learning at each stage?

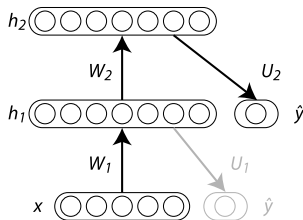
- Each layer is trained as the hidden layer of a supervised 2-layer neural net.
- After training the 2-layer neural net, discard output layer;
- Propagate data through new hidden layer, train another layer, etc.



Greedy Supervised Layer-Wise Initialization

Why not use supervised learning at each stage?

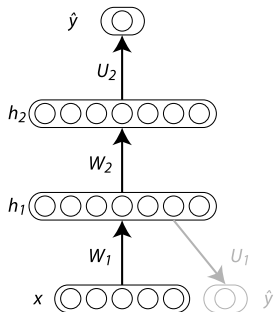
- Each layer is trained as the hidden layer of a supervised 2-layer neural net.
- After training the 2-layer neural net, discard output layer;
- Propagate data through new hidden layer, train another layer, etc.



Greedy Supervised Layer-Wise Initialization

Why not use supervised learning at each stage?

- Each layer is trained as the hidden layer of a supervised 2-layer neural net.
- After training the 2-layer neural net, discard output layer;
- Propagate data through new hidden layer, train another layer, etc.



Experiments on Greedy Layer-Wise Initialization

	train.	test.
Deep Belief Net, unsupervised pre-training	0%	1.2%
Deep net, autoassociator pre-training	0%	1.4%
Deep net, supervised pre-training	0%	2.0%
Deep net, no pre-training	.004%	2.4%
Shallow net, no pre-training	.004%	1.9%

Classification error on MNIST digits benchmark training, validation, and test sets, with the best hyper-parameters according to validation error.

Deep nets with 3 to 5 hidden layers.
Selects around 500 hidden units per layer.

Experiments on Greedy Layer-Wise Initialization

	train.	test.
Deep Belief Net, unsupervised pre-training	0%	1.2%
Deep net, autoassociator pre-training	0%	1.4%
Deep net, supervised pre-training	0%	2.0%
Deep net, no pre-training	.004%	2.4%
Shallow net, no pre-training	.004%	1.9%

Classification error on MNIST digits benchmark training, validation, and test sets, with the best hyper-parameters according to validation error.

Deep nets with 3 to 5 hidden layers.
Selects around 500 hidden units per layer.

Supervised greedy is **too greedy**.
Greedy unsupervised initialization works great.

Is it Really an Optimization Problem?

Why 0 train error even with deep net / no-pretraining?

Is it Really an Optimization Problem?

Why 0 train error even with deep net / no-pretraining?

Classification error on MNIST with 20 hidden units on top layer:

	train.	test.
Deep net, autoassociator pre-training	0%	1.6%
Deep net, supervised pre-training	0%	1.9%
Deep net, no pre-training	.59%	2.2%
Shallow net, no pre-training	3.6%	5.0%

Is it Really an Optimization Problem?

Why 0 train error even with deep net / no-pretraining?

Classification error on MNIST with 20 hidden units on top layer:

	train.	test.
Deep net, autoassociator pre-training	0%	1.6%
Deep net, supervised pre-training	0%	1.9%
Deep net, no pre-training	.59%	2.2%
Shallow net, no pre-training	3.6%	5.0%

Because

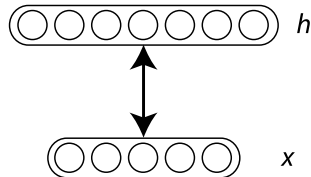
- 1 last fat hidden layer did all the work
- 2 using a poor representation (output of all previous layers)

Yes it is really an **optimization** problem
and a **representation** problem

Restricted Boltzmann Machines

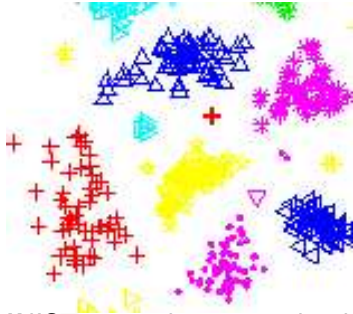
Bi-partite Boltzmann machine:

$$P(V = v, H = h) \propto e^{-\mathcal{E}(v,h)} = e^{v'b+h'c+h'Wv}$$



- Conditionals $P(v|h)$ and $P(h|v)$ easy to derive, and factorize.
- Contrastive divergence provides good estimator of log-likelihood gradient.
- Originally for binary variables; we extend it **easily** to continuous variables by slightly changing energy function and range of values (see poster).

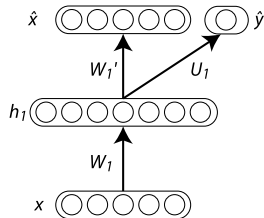
Combining Supervised & Unsupervised



MNIST: nice clusters in the distribution

⇒ **input distribution structure reveals the target class.**

Combining Supervised & Unsupervised



MNIST: nice clusters in the distribution

⇒ **input distribution structure reveals the target class.**

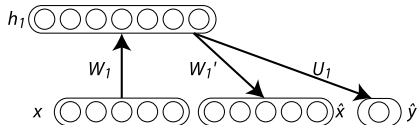
Otherwise? Simple solution:

combine supervised & unsupervised layer-wise greedy initialization.

Just add the two stochastic gradient updates.

Similar to Partial Least Squares

Combining Supervised & Unsupervised



MNIST: nice clusters in the distribution

⇒ **input distribution structure reveals the target class.**

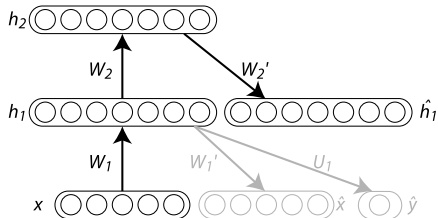
Otherwise? Simple solution:

combine supervised & unsupervised layer-wise greedy initialization.

Just add the two stochastic gradient updates.

Similar to Partial Least Squares

Combining Supervised & Unsupervised



MNIST: nice clusters in the distribution

⇒ **input distribution structure reveals the target class.**

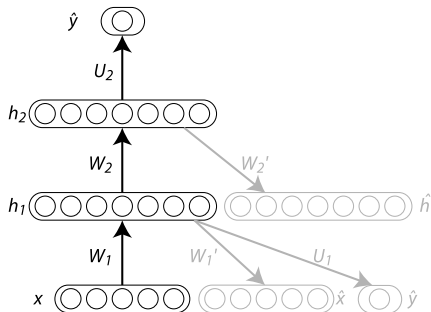
Otherwise? Simple solution:

combine supervised & unsupervised layer-wise greedy initialization.

Just add the two stochastic gradient updates.

Similar to Partial Least Squares

Combining Supervised & Unsupervised



MNIST: nice clusters in the distribution

⇒ **input distribution structure reveals the target class.**

Otherwise? Simple solution:

combine supervised & unsupervised layer-wise greedy initialization.

Just add the two stochastic gradient updates.

Similar to Partial Least Squares

More experimental results

	Abalone MSE	Cotton class. error
1. Deep Network, no pre-training	4.2	43.0%
2. Logistic regression	.	45.0%
3. DBN, Bin inputs, unsup	4.47	45.0%
4. DBN, Bin inputs, partially sup	4.28	43.7%
5. DBN, Gauss inputs, unsup	4.19	35.8%
6. DBN, Gauss inputs, partially sup	4.18	31.4%

MSE on Abalone task and classification error on Cotton task, showing improvement with Gaussian vs binomial units and partial supervision

Conclusions

- For AI \Rightarrow must learn **highly-varying functions** efficiently
 \Rightarrow **deep architectures** (statistical efficiency)

Conclusions

- For AI \Rightarrow must learn **highly-varying functions** efficiently
 \Rightarrow **deep architectures** (statistical efficiency)
- Deep architectures not trainable? computational efficiency?
new methods appear to **break through the obstacle**

Conclusions

- For AI \Rightarrow must learn **highly-varying functions** efficiently
 \Rightarrow **deep architectures** (statistical efficiency)
- Deep architectures not trainable? computational efficiency?
new methods appear to **break through the obstacle**
- Basic principle: **greedy layer-wise unsupervised** (or adding unsupervised and supervised criteria)

Conclusions

- For AI \Rightarrow must learn **highly-varying functions** efficiently
 \Rightarrow **deep architectures** (statistical efficiency)
- Deep architectures not trainable? computational efficiency?
new methods appear to **break through the obstacle**
- Basic principle: **greedy layer-wise unsupervised** (or adding unsupervised and supervised criteria)
- **Principle works about as well with symmetric autoassociators** in feedforward neural net

Conclusions

- For AI \Rightarrow must learn **highly-varying functions** efficiently
 \Rightarrow **deep architectures** (statistical efficiency)
- Deep architectures not trainable? computational efficiency?
new methods appear to **break through the obstacle**
- Basic principle: **greedy layer-wise unsupervised** (or adding unsupervised and supervised criteria)
- **Principle works about as well with symmetric autoassociators** in feedforward neural net
- The **unsupervised part** is important: regularizes and makes sure to propagate most information about input, **purely supervised is too greedy**.

Conclusions

- For AI \Rightarrow must learn **highly-varying functions** efficiently
 \Rightarrow **deep architectures** (statistical efficiency)
- Deep architectures not trainable? computational efficiency?
new methods appear to **break through the obstacle**
- Basic principle: **greedy layer-wise unsupervised** (or adding unsupervised and supervised criteria)
- **Principle works about as well with symmetric autoassociators** in feedforward neural net
- The **unsupervised part** is important: regularizes and makes sure to propagate most information about input, **purely supervised is too greedy**.
- Easy extensions of Deep Belief Nets: continuous-valued units / partially supervised initialization when input density is not revealing of target