Learning Deep Architectures for AI

Yoshua Bengio

May 2008

Thanks to : James Bergstra, Olivier Breuleux, Aaron Courville, Olivier Delalleau, Dumitru Erhan, Pascal Lamblin, Hugo Larochelle, Jerome Louradour, Nicolas Le Roux, Pierre-Antoine Manzagol, Dan Popovici, François Rivest, Joseph Turian, Pascal Vincent Check review paper (with same title) on my web page

Artificial Intelligence

- Half century of research and goal seems still so far
- Why?
- Too much in a hurry for results rather than understanding?



Knowledge : from Where?

- Al needs much knowledge about our world
- Explicit-symbolic approach :

 $\ensuremath{\text{Cyc}}$ = hand-crafted collection of rules and facts



Gigantic



Incoherent/Incomplete



Not robust not uncertainty-friendly

3

Learning the Knowledge?

- Animals and humans : innate + learned
- Can learn many tasks that were not evolution-tuned !
- Bet on existence of some generic strategies/principles
- More exciting / greater payoff / worth exploring this path



Compact Representation \Rightarrow Generalization

Extract the essence from the data \Rightarrow generalization

• Occam's Razor

- Kolmogorov/Solomonoff 1964
- Vapnik 1972





 $\min_{f} \frac{1}{n} \sum_{i=1}^{n} error(x_i, f(x_i)) \neq \\ \min_{f} E[error(x_i, f(x_i))]$

1D Generalization

Easy if the target function is smooth (few variations).



< 回 > < 三 > < 三 >

Works well with a good representation : where notion of neighborhood is meaningful.



- 4 同 6 4 日 6 4 日 6

Where Local Generalization Fails



- Pixel-to-pixel Euclidean distance only works very locally
- Klingons work in a more abstract representation space

(人間) システン イラン

Local Generalization : negative results

We and others have shown negative results illustrating limitations of local generalization :

- (classical non-parametric)
- Local kernel methods :
 - SVMs
 - Gaussian Process
 - Graph-based semi- and un-supervised learning alg. (LLE, Isomap, etc.)
- Decision trees

All break data space in regions s.t. # degrees freedom $\propto \#$ regions





伺 ト イヨト イヨト

Curse of Dimensionality

Basic result :

Theorem Sketch

Local learning algorithms cannot generalize to variations not covered by the training set.

Informal Corollary

Local learning algorithms can require a number of training examples exponential in the input dimension to obtain a given generalization error.



(日) (同) (三) (三)

Actual theoretical results (NIPS'2004, NIPS'2005, *Semi-Supervised Learning* book) are specialized :

- Gaussian kernel machines
 - Functions varying a lot along a straight line
 - Parity function
- Semi-supervised learning from neighborhood-graph
- Decision trees on highly varying functions
- Local kernel manifold learning from neighborhood-graph (such as kPCA, LLE, Isomap, ...)

イロト イポト イヨト イヨト 二三

Primates Visual System



Visual System : sequence of transformations / abstraction levels

・ロン ・回 と ・ヨン ・ヨン

Strategy : Distributed Representations

- Distributed representation : each percept represented by the combination of many features
- Multiple levels of representation (like in the brain)
- Exponential advantage
- Missing in many current learning algorithms (most clustering, non-parametric, semi-supervised, kernels, mixture models, etc.)



< 回 > < 三 > < 三 >

Exploiting Multiple Levels of Representation



Higher-level abstractions

Primitive pattern detectors

Oriented edge detectors

伺 ト イヨト イヨト

3

Pixels

Computation Graph and Depth

Each node = computation element (from a set)



compute
$$x * sin(a * x + b)$$
,
depth=4.

elements = artificial neurons f(x) = tanh(b + w'x). Multi-layer neural net depth=3.

Theorem Sketch

When a function can be compactly represented by a deep architecture, representing it with a shallow architecture can require a number of elements exponential in the input dimension.



- 4 同 6 4 日 6 4 日 6

- 2-level logic circuits on $\{0,1\}^n$
 - can represent any discrete function
 - most functions require $O(2^n)$ logic gates
 - ∃ functions computable efficiently with depth k, requiring O(2ⁿ) gates if depth ≤ k − 1.



- < 同 > < 三 > < 三 >

• Similar results for circuits of formal neurons

- ∃ "simple" functions requiring exponential size architectures with Gaussian kernel machines
- Most current/popular learning algorithms : depth = 1 ou 2, sometimes 3 (e.g. boosted trees, forests, 2-hidden layer MLPs)



Insufficient Depth : Consequence

- Theoretical results (proved for some element sets) insufficient depth ⇒ very fat architecture
- ullet \Rightarrow very large number of parameters required
- ⇒ very large number of examples required



🗇 🕨 🖌 🚍 🕨 🤇 🚍 🕨

Training Deep Architectures : the Challenge

- Two levels suffice to represent any function
- Up to 2006, failure of attempts to train deep architectures
- Why? Non-convex optimisation and stochastic!
- Focus NIPS 1995-2005 : convex learning algorithms





 \Rightarrow Let us face the challenge !

Strategy : multi-task semi-supervised learning

- Humans : most examples not semantically labeled
- Each example informs many tasks
- Representations shared among tasks
- Capture common factors of variations to generalize easily to new tasks (even with 0 examples !)
- $\lim_{\# tasks \to \infty} = (un+semi)$ -supervised



Strategy : one level of abstraction at a time

- Humans learn simple things first : first levels of visual system converge in critical periods
- Percept representation = abstraction of the percept
- Learn a first level of representation, then a second built on the first, etc.



2006 : Breakthrough !

- FIRST : successful training of deep architectures ! Hinton et al Neural Comp. 2006, followed by Bengio et al, and Ranzato et al at NIPS'2006
- One trains one layer after the other of an MLP
- Unsupervised learning in each layer of initial representation
- Continue training an ordinary but deep MLP near a better minimum



Individual Layer : RBMs and auto-encoders





◆□▶ ◆□▶ ◆三▶ ◆三▶ 三 のので

Learning deep networks Supervised fine-tuning

- Initial deep hierarchical mapping is learnt in an unsupervised way.
- \rightarrow initialization for a supervised task.
- Output layer gets added.
- Global fine tuning by gradient descent on supervised criterion.



- 4 回 ト - 4 回 ト

Learning deep networks Supervised fine-tuning

- Initial deep hierarchical mapping is learnt in an unsupervised way.
- \rightarrow initialization for a supervised task.
- Output layer gets added.
- Global fine tuning by gradient descent on supervised criterion.



(人間) システン イラン

Learning deep networks Supervised fine-tuning

- Initial deep hierarchical mapping is learnt in an unsupervised way.
- \rightarrow initialization for a supervised task.
- Output layer gets added.
- Global fine tuning by gradient descent on supervised criterion.



< 回 > < 三 > < 三 >

Deep Belief Nets

- Hinton et al. (2006) introduced the Deep Belief Network (DBN), a deep probabilistic/generative neural network
- The training procedure is first **layer-wise greedy** and **unsupervised** (initialization).



 Then the model is converted into a conditional predictor and fine-tuned

$$\min_{\theta} -\frac{1}{n} \sum_{i=t}^{n} \log \hat{p}(y_t | x_t, \theta)$$



- Clean input x ∈ [0, 1]^d is partially destroyed, yielding corrupted input : x̃ ~ q_D(x̃|x).
- $\tilde{\mathbf{x}}$ is mapped to hidden representation $\mathbf{y} = f_{\theta}(\tilde{\mathbf{x}})$.
- From **y** we reconstruct a $\mathbf{z} = g_{\theta'}(\mathbf{y})$.
- Train parameters to minimize the cross-entropy "reconstruction error"
- Corresponds to maximizing variational bound on likelihood of a generative model
- Naturally handles missing values / occlusion / multi-modality

(ロ) (部) (E) (E) (E)



- Clean input x ∈ [0, 1]^d is partially destroyed, yielding corrupted input : x̃ ~ q_D(x̃|x).
- $\tilde{\mathbf{x}}$ is mapped to hidden representation $\mathbf{y} = f_{\theta}(\tilde{\mathbf{x}})$.
- From **y** we reconstruct a $\mathbf{z} = g_{\theta'}(\mathbf{y})$.
- Train parameters to minimize the cross-entropy "reconstruction error"
- Corresponds to maximizing variational bound on likelihood of a generative model

(ロ) (部) (E) (E) (E)

• Naturally handles missing values / occlusion / multi-modality

Denoising Auto-Encoders



- Clean input x ∈ [0, 1]^d is partially destroyed, yielding corrupted input : x̃ ~ q_D(x̃|x).
- $\tilde{\mathbf{x}}$ is mapped to hidden representation $\mathbf{y} = f_{\theta}(\tilde{\mathbf{x}})$.
- From **y** we reconstruct a $\mathbf{z} = g_{\theta'}(\mathbf{y})$.
- Train parameters to minimize the cross-entropy "reconstruction error"
- Corresponds to maximizing variational bound on likelihood of a generative model

・ロト ・回ト ・ヨト ・ヨト

• Naturally handles missing values / occlusion / multi-modality

Denoising Auto-Encoders



- Clean input x ∈ [0, 1]^d is partially destroyed, yielding corrupted input : x̃ ~ q_D(x̃|x).
- $\tilde{\mathbf{x}}$ is mapped to hidden representation $\mathbf{y} = f_{\theta}(\tilde{\mathbf{x}})$.
- From **y** we reconstruct a $\mathbf{z} = g_{\theta'}(\mathbf{y})$.
- Train parameters to minimize the cross-entropy "reconstruction error"
- Corresponds to maximizing variational bound on likelihood of a generative model
- Naturally handles missing values / occlusion / multi-modality

・ロン ・回 と ・ヨン ・ヨン

Denoising Auto-Encoders



- Clean input x ∈ [0, 1]^d is partially destroyed, yielding corrupted input : x̃ ~ q_D(x̃|x).
- $\tilde{\mathbf{x}}$ is mapped to hidden representation $\mathbf{y} = f_{\theta}(\tilde{\mathbf{x}})$.
- From **y** we reconstruct a $\mathbf{z} = g_{\theta'}(\mathbf{y})$.
- Train parameters to minimize the cross-entropy "reconstruction error"
- Corresponds to maximizing variational bound on likelihood of a generative model

イロト 不得 トイヨト イヨト 二日

• Naturally handles missing values / occlusion / multi-modality

Manifold Learning Perspective



Denoising autoencoder can be seen as a way to learn a manifold :

- Suppose training data (x) concentrate near a low-dimensional manifold.
- Corrupted examples (•) are obtained by applying corruption process $q_{\mathcal{D}}(\widetilde{X}|X)$ and will lie farther from the manifold.
- The model learns with $p(X|\tilde{X})$ to "project them back" onto the manifold.
- Intermediate representation Y can be interpreted as a coordinate system for points on the manifold.

<ロ> <同> <同> < 回> < 回>

basic : subset of original MNIST digits : 10 000 training samples, 2 000 validation samples, 50 000 test samples.



(a) **rot** : applied random rotation (angle between 0 and 2π radians)



(c) **bg-img** : background is random patch from one of 20 images



(b) **bg-rand :** background made of random pixels (value in $0 \dots 255$)



(d) **rot-bg-img :** combination of rotation and background image

Benchmark problems

• rect : discriminate between tall and wide rectangles on black background.



- rect-img : borderless rectangle filled with random image patch. Background is a different image patch.
- convex : discriminate between convex and non-convex shapes.



Dataset	SVM _{rbf}	DBN-3	SAA-3	SdA-3 (ν)
basic	3.03±0.15	$3.11{\scriptstyle \pm 0.15}$	$3.46{\scriptstyle \pm 0.16}$	$2.80_{\pm 0.14} \ (10\%)$
rot	$11.11_{\pm 0.28}$	10.30±0.27	10.30±0.27	$10.29_{\pm 0.27}$ (10%)
bg-rand	$14.58 \scriptstyle \pm 0.31$	6.73±0.22	$11.28{\scriptstyle \pm 0.28}$	$10.38_{\pm 0.27}$ (40%)
bg-img	22.61±0.37	$16.31{\scriptstyle\pm0.32}$	23.00 _{±0.37}	16.68±0.33 (25%)
rot-bg-img	55.18 _{±0.44}	$47.39_{\pm 0.44}$	51.93 _{±0.44}	44.49 ±0.44 (25%)
rect	2.15±0.13	$2.60{\scriptstyle \pm 0.14}$	$2.41{\scriptstyle \pm 0.13}$	$1.99_{\pm 0.12}$ (10%)
rect-img	24.04 _{±0.37}	22.50 _{±0.37}	24.05 _{±0.37}	21.59 _{±0.36} (25%)
convex	19.13 _{±0.34}	18.63±0.34	$18.41{\scriptstyle \pm 0.34}$	$19.06_{\pm 0.34} \ (10\%)$

Some Results 2006-2008

- Deep architectures superior on MNIST (NIPS'2006)
- Greater advantage on more complex tasks (ICML'2007)
- RBMs slightly better than ordinary auto-encoders
- \bullet Denoising auto-encoders \geq RBMs, and more flexible
- Applications in NLP, vision, MOCAP, collaborative filtering
- Optimization sometimes still deficient, challenges ahead





◆□ > ◆□ > ◆ 三 > ◆ 三 > ● の < @ >

Several Strategies are Continuations

- Older : stochastic gradient from small parameters
- Breakthrough : greedy layer-wise construction
- New : gradually bring in more difficult examples



- 4 回 ト - 4 回 ト

Curriculum Strategy

Start with simpler, easier examples, and gradually introduce more of the more complicated ones as the learner is ready to learn them.



Design the sequence of tasks / datasets to guide learning/optimization.



Strategy : Society = Parallel Optimisation

- Each human = potential solution
- Better solutions spread through language
- Similar to genetic evolution : parallel search + recombination
- R. Dawkins' Memes
- Simulations support this hypothesis
- AI : take avantage of human culture



Combine many strategies, to obtain a baby AI that masters the semantics of a simple visual + linguistic universe



Subject	Question	Answer		
Color	There is a small triangle. What color is it?	Green		
Shape	What is the shape of the green object?	Triangle		
Location	Is the blue square at the top or at the bottom?	At the top		
Size	There is a triangle on the right.			
	ls it rather small or big?	Small		
Size (relative)	Is the square smaller or bigger than the triangle?	Bigger		

- 4 回 ト - 4 回 ト

Conclusions and work in progress

- Al \Rightarrow learning \Rightarrow generalize
 - \Rightarrow generalize non-locally \Rightarrow learn distributed representations
 - \Rightarrow deep architectures \Rightarrow optimisation challenge
- Breakthrough in 2006
- Biological inspiration : humans' strategies to optimize learning of world model
 - multi-task, unsupervised, semi-supervised
 - multiple levels of distributed representation
 - learn lower levels first
 - curriculum / education
 - collective parallel search
- Ongoing work : denoising auto-encoders, work as well or better than DBNs
- Collobert & Weston : learning representations by layer-wise manifold learning

• Patience : see the long term...