



Yoshua Bengio&Olivier DelalleauStatistical Learning Algorithms Canada Research Chair, U. Montreal

ALT / DS 2011 October 5th, 2011, Espoo, Finland

ON THE EXPRESSIVE POWER OF DEEP ARCHITECTURES

From AI to Learning

- Intelligence requires knowledge \rightarrow decisions
- Knowledge can be implicit
- Explicitly providing knowledge failed (expert systems)
 - Verbally expressed knowledge is incomplete
 - And lacks the required expression of uncertainty
- Learning captures knowledge from data
 - Can capture what is needed (completeness)
 - Can be actionable (learn to predict & act)
 - Can handle uncertainty (probabilistic models)

Learning to Generalize. How?

- Capturing dependencies between random variables
- Spreading out the probability mass from the empirical distribution. Where???
- Discovering underlying abstractions / explanatory factors

Onto Deep Learning

- Real-world distributions have convoluted unknown structure, not all captured by the principle of local generalization
- We want weak priors that are stronger than the usual smoothness prior
- Deep Learning: a way to address this by the discovery of multiple levels of representation capturing the underlying factors of variation

Shallow learning architecture



1-layer NNet, SVM, GP predictor, decision tree, boosted stumps, etc.

Deep learning architecture



Deep Motivations

- Brains have a deep architecture
- Cortex seems to have a generic learning algorithm
- Humans' ideas composed from simpler ones
- Insufficient depth can be exponentially inefficient
- Distributed (possibly sparse) representations necessary for nonlocal generalization, exponentially more efficient than 1-of-N enumeration of latent variable values
- Multiple levels of latent variables allow combinatorial sharing of statistical strength





Deep Architecture in our Mind

- Humans organize their ideas and concepts hierarchically
- Humans first learn simpler concepts and then compose them to represent more abstract ones
- Engineers break-up solutions into multiple levels of abstraction and processing
- It would be nice to learn / discover these concepts

(knowledge engineering failed because of limits of introspection?)



Deep Learning Hypotheses

Hypothesis 1: deep hierarchy of features useful to efficiently represent and learn complex abstractions needed for AI and mammal intelligence.

- Computational & statistical efficiency •
- allexamples Hypothesis 2: unsupervised learning of representations is Even if we had labels for a crucial component of the solution.
 - **Optimization & regularization.** •
- Theoretical and ML-experimental support for both.

Principle of Local Generalization



The Curse of Dimensionality

To generalize locally, need representative examples for all relevant variations!

Classical solution: hope for a smooth enough target function



Limits of Local Generalization: Theoretical Results





e.g. Gaussian (RBF) SVM

 Theorem: Gaussian kernel machines need at least k examples to learn a function that has 2k zerocrossings along some line



 Theorem: For a Gaussian kernel machine to learn some maximally varying functions over *d* inputs requires O(2^d) examples

Curse of Dimensionality When Generalizing Locally on a Manifold (Bengio et al 2006)



How to Beat the Curse of Many Factors of Variation?

Compositionality: exponential gain in representational power

- Distributed representations / embeddings: feature learning
- Deep architecture: multiple levels of feature learning

Can generalize to new configurations

Distributed Representations

- Many features active simultaneously
- Input represented by the activation of a set of features that are not mutually exclusive
- Can be exponentially more efficient than local representations
- FEATURE LEARNING instead of / on top of manual feature-engineering

Local vs Distributed Latent Variables

Exponentially more regions can be distinguished for the same number of parameters, i.e., examples



RBM Hidden Units Carve Input Space



Restricted Boltzmann Machine

 The most popular building block for deep architectures

$$P(x,h) = \frac{1}{Z}e^{b^T h + c^T x + h^T W x}$$

- Bipartite undirected graphical model
- Inference is trivial:
- P(h|x) & P(x|h) factorize



Discrete Input RBMs are Universal Approximators



- Adding one hidden unit (with proper choice of parameters) guarantees increasing likelihood
- With enough hidden units, can perfectly model any discrete distribution
- RBMs with variable nb of hidden units = non-parametric

Continuous Inputs - The Best Generative Model of Images: Spike-and-Slab RBM

Samples from μ -ssRBM:



Nearest examples in CIFAR: (least square dist.)



2011





subroutine1 includes subsub1 code and subsub2 code and subsubsub1 code

subroutine2 includes subsub2 code and subsub3 code and subsubsub3 code and ...



"Shallow" computer program

Architecture Depth



"Deep" circuit



"Shallow" circuit



Falsely reassuring theorems: one can approximate any reasonable (smooth, boolean, etc.) function with a 2-layer architecture

Deep Architectures are More Expressive



Sharing Components in a Deep Architecture Polynomial expressed with shared components: advantage of depth may grow exponentially





. Depth 2 suffices to represent any finite polynomial (sum of products)

. (Poon & Domingos 2010) use deep sum-product networks to efficiently parametrize partition functions

Polynomials that Need Depth



* Need O(n) nodes with depth log(n) circuit

- 2i layers and $n = 4^i$ input variables
- alternate additive and multiplicative units • unit ℓ_i^k takes as inputs ℓ_{2i-1}^{k-1} and ℓ_{2i}^{k-1}

* Need O($2^{\sqrt{n}}$) nodes with depth-2 circuit

More Polynomials that Need Depth



- 2i + 1 layers and n variables (n independent of i)
- alternate multiplicative and additive units
- unit ℓ_j^k takes as inputs $\{\ell_m^{k-1} | m \neq j\}$

More Deep Theory

Poly-logarithmic Independence Fools Bounded-Depth Boolean Circuits Braverman, CACM 54(4), April 2011.

If all marginals of the input distribution involving at most k variables are uniform, higher depth makes it exponentially easier to distinguish the joint from the uniform.

Deep Architectures and Sharing Statistical Strength, Multi-Task Learning

- Generalizing better to new tasks is crucial to approach AI
- Deep architectures learn (good intermediate representations that can be shared across tasks
- Good representations make sense for many tasks





Parts Are Re-Used to Form Different Objects

Layer 3: objects

Layer 2: parts

Layer 1: edges

(Lee et al. ICML 2009)

Before 2006

Failing to train deep architectures



2006: The Deep Breakthrough



- Hinton, Osindero & Teh « <u>A Fast Learning</u> <u>Algorithm for Deep</u> <u>Belief Nets</u> », Neural Computation, 2006
- Bengio, Lamblin, Popovici, Larochelle « <u>Greedy Layer-Wise</u> <u>Training of Deep</u> <u>Networks</u> », *NIPS'2006*
 - Ranzato, Poultney, Chopra, LeCun « Efficient Learning of Sparse Representations with an Energy-Based Model », NIPS'2006

Deep training

input OOO ... O















Supervised Fine-Tuning



Stacking Auto-Encoders



Palette of Tricks to Train Energy-Based Models

Partition function expensive (vocab.) or intractable

- Contrastive Divergence
- PCD / SML + MCMC tricks
 - Tempering
 - Mean-field / variational, etc.
- (regularized) Score Matching / denoising
- Sparse coding / Sparse Predictive Decomposition
- Ratio Matching
- Pseudo-likelihood
- Ranking / margin-based criteria
- Noise contrastive estimation

Most rely on + vs – examples contrast

See my book / review paper (F&TML 2009): Learning Deep Architectures for AI

Sparse Auto-Encoders & Sparse Coding

- Penalty on the representation to achieve sparsity.
- Stacked sparse auto-encoders successfully used by Andrew Ng's group at Stanford (e.g. ICML 2011)
- Used by Google in their Google Goggles vision system
- Sparse coding (recently stacked as well)
- Sparse Predictive Decomposition (LeCun)

Denoising Auto-Encoder

(Vincent et al 2008, 2010)



- Stochastically corrupt the input
- Reconstruction target = clean input



Stacked Denoising Auto-Encoders

- No partition function, can measure training criterion
- Encoder & decoder: any parametrization
- As good or better than RBMs for feature learning
- = regularized score matching



Unsupervised and Transfer Learning Challenge: 1st Place in Final Phase



Contractive Auto-Encoders



cannot afford contraction in manifold directions Training criterion:

wants contraction in all directions

Few active units

$$\mathcal{J}_{CAE}(\theta) = \sum_{x \in D_n} \left(L(x, g(h(x))) + \lambda \sum_{ij} \left(\frac{\partial h_j(x)}{\partial x_i} \right)^2 \right)$$

Manifold Tangent Classifier (NIPS 2011)

• Leading singular vectors on MNIST, CIFAR-10, RCV1:



Knowledge-free MNIST: 0.81% error

K-NN	NN	SVM	DBN	CAE	DBM	CNN	MTC
3.09%	1.60%	1.40%	1.17%	1.04%	0.95%	0.95%	0.81%

Unsupervised Learning: Disentangling Factors of Variation

- (Goodfellow et al NIPS'2009): some hidden units more invariant (with more depth) to input geometry variations
- (Glorot et al ICML'2011): some hidden units specialize on one aspect (domain) while others on another (sentiment)
- We don't want invariant representations because it is not clear to what aspects, but disentangling factors would help a lot
- Sparse/saturated units seem to help
- Why?
- How to train more towards that objective?

Recent Deep Learning Highlights

- Google Goggles uses stacked sparse autoencoders (Hartmut Neven @ ICML 2011)
- UofT breaks old accuracy ceiling in TIMIT phoneme detection
- Stanford breaks records in video / gesture classification
- NYU breaks records in traffic sign class
- Montreal wins Unsupervised & Transfer Learning Challenge

Conclusions

- Deep Learning: powerful arguments & generalization principles
- Unsupervised Feature Learning is crucial: many new algorithms and applications in recent years
- DL particularly suited for multi-task learning, transfer learning, domain adaptation, self-taught learning, and semisupervised learning with few labels

http://deeplearning.net

HOME TUTORIALS ABOUT READING LIST SOFTWARE LINKS BLOG DEMOS DATASETS EVENTS BIBLIOGRAPHY



Deep Learning

... moving beyond shallow machine learning since 2006!



LISA Lab Wins the Final Phase of UTLC Challenge New Challenge Announced Deep Learning Workshop at NIPS 2010 New Events Page Deep Learning papers at ICML 2010

Tags

Meta

Log in Entries RSS Comments RSS WordPress.org

Welcome to Deep Learning

Deep Learning is a new area of Machine Learning research, which has been introduced with the objective of moving Machine Learning closer to one of its original goals: Artificial Intelligence.

This website is intended to host a variety of resources and pointers to information about Deep Learning. In these pages you will find

- a reading list,
- · links to software,
- datasets,
- a discussion forum,
- as well as tutorials and cool demos.

For the latest additions, including papers and software announcement, **be sure to visit the Blog section** of the website. **Contact us** if you have any comments or suggestions!

Last modified on April 26, 2010, at 9:45 am by ranzato

Pages

About Bibliography Blog Datasets Demos Events Reading List Software links Tutorials

Links

Discussion forum

<u>http://deeplearning.net/software/theano</u> : numpy \rightarrow GPU

