

An Empirical Evaluation of Deep Architectures on Problems with Many Factors of Variation

Hugo Larochelle, Dumitru Erhan, Aaron Courville,
James Bergstra, and Yoshua Bengio
Université de Montréal

13/06/2007

- Started in 2005, for a total of 58 talks + 2 “Best-of NIPS” sessions.
- You can find most of the past talk slides and information about upcoming talks on the web page
- Bilingual page :
 - “**McGill-UdeM-MITACS** Machine Learning Seminars” in english
 - “Séminaires **UdeM-McGill-MITACS** d’Apprentissage Automatique” in french
- Google keywords :
 - “mitacs machine learning”
 - “active learning in POMDP”
 - “echo state network apprentissage”
 - “Aaron Courville”
 - “aaron courville site :umontreal.ca”
 - “indian sumit chopra”
 - “an efficient method for gradient based applications SVM models”

- Started in 2005, for a total of 58 talks + 2 “Best-of NIPS” sessions.
- You can find most of the past talk slides and information about upcoming talks on the web page
- Bilingual page :
 - “**McGill-UdeM-MITACS** Machine Learning Seminars” in english
 - “Séminaires **UdeM-McGill-MITACS** d’Apprentissage Automatique” in french
- Google keywords :
 - “mitacs machine learning”
 - “active learning in POMDP”
 - “echo state network apprentissage”
 - “Aaron Courville”
 - “aaron courville site :umontreal.ca”
 - “indian sumit chopra”
 - “an efficient method for gradient based applications SVM models”

- Started in 2005, for a total of 58 talks + 2 “Best-of NIPS” sessions.
- You can find most of the past talk slides and information about upcoming talks on the web page
- Bilingual page :
 - “**McGill-UdeM-MITACS** Machine Learning Seminars” in english
 - “Séminaires **UdeM-McGill-MITACS** d’Apprentissage Automatique” in french
- Google keywords :
 - “mitacs machine learning”
 - “active learning in POMDP”
 - “echo state network apprentissage”
 - “Aaron Courville”
 - “aaron courville site :umontreal.ca”
 - “indian sumit chopra”
 - “an efficient method for gradient based applications SVM models”

- Started in 2005, for a total of 58 talks + 2 “Best-of NIPS” sessions.
- You can find most of the past talk slides and information about upcoming talks on the web page
- Bilingual page :
 - “**McGill-UdeM-MITACS** Machine Learning Seminars” in english
 - “Séminaires **UdeM-McGill-MITACS** d’Apprentissage Automatique” in french
- Google keywords :
 - “mitacs machine learning”
 - “active learning in POMDP”
 - “echo state network apprentissage”
 - “Aaron Courville”
 - “aaron courville site :umontreal.ca”
 - “indian sumit chopra”
 - “an efficient method for gradient based applications SVM models”

Introduction

- Recently, models relying on deep architectures have been proposed (Deep Belief Networks)
- Their performance compare favorably to state of art models such as Support Vector Machines
- They have been tested on relatively few and simple problems
- We propose to evaluate them on problems with many factors of variation

Introduction

- Recently, models relying on deep architectures have been proposed (Deep Belief Networks)
- Their performance compare favorably to state of art models such as Support Vector Machines
- They have been tested on relatively few and simple problems
- We propose to evaluate them on problems with many factors of variation

Introduction

- Recently, models relying on deep architectures have been proposed (Deep Belief Networks)
- Their performance compare favorably to state of art models such as Support Vector Machines
- They have been tested on relatively few and simple problems
- We propose to evaluate them on problems with many factors of variation

Introduction

- Recently, models relying on deep architectures have been proposed (Deep Belief Networks)
- Their performance compare favorably to state of art models such as Support Vector Machines
- They have been tested on relatively few and simple problems
- We propose to evaluate them on problems with many factors of variation

Problems with Many Factors of Variation

- We focus here on problems where the input distribution has the following structure

$$p(x) = \sum_{\phi_1, \dots, \phi_m} p(x|\phi_1, \dots, \phi_m)p(\phi_1, \dots, \phi_m)$$

where $p(\phi_1, \dots, \phi_m)$ is high for (exponentially) many combinations of values of the factors of variation ϕ_i

- Problems with such a structure :
 - digit recognition (vision) : $\phi_i \in \{\text{rotation angle, scaling, background, etc.}\}$
 - document classification (NLP) : $\phi_i \in \{\text{topics, style, etc.}\}$
 - speech recognition (signal processing) :
 $\phi_i \in \{\text{speaker gender, background noise, environment echo, etc.}\}$
- We will focus on vision problems
- We want to avoid hand-engineered solutions to these problems

Problems with Many Factors of Variation

- We focus here on problems where the input distribution has the following structure

$$p(x) = \sum_{\phi_1, \dots, \phi_m} p(x|\phi_1, \dots, \phi_m)p(\phi_1, \dots, \phi_m)$$

where $p(\phi_1, \dots, \phi_m)$ is high for (exponentially) many combinations of values of the factors of variation ϕ_i

- Problems with such a structure :
 - digit recognition (vision) : $\phi_i \in \{\text{rotation angle, scaling, background, etc.}\}$
 - document classification (NLP) : $\phi_i \in \{\text{topics, style, etc.}\}$
 - speech recognition (signal processing) :
 $\phi_i \in \{\text{speaker gender, background noise, environment echo, etc.}\}$
- We will focus on vision problems
- We want to avoid hand-engineered solutions to these problems

Problems with Many Factors of Variation

- We focus here on problems where the input distribution has the following structure

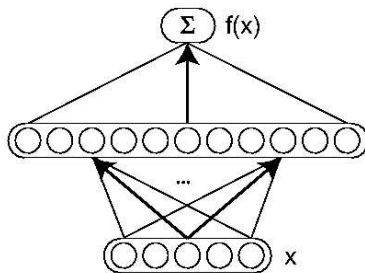
$$p(x) = \sum_{\phi_1, \dots, \phi_m} p(x|\phi_1, \dots, \phi_m)p(\phi_1, \dots, \phi_m)$$

where $p(\phi_1, \dots, \phi_m)$ is high for (exponentially) many combinations of values of the factors of variation ϕ_i

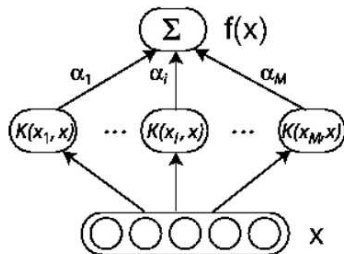
- Problems with such a structure :
 - digit recognition (vision) : $\phi_i \in \{\text{rotation angle, scaling, background, etc.}\}$
 - document classification (NLP) : $\phi_i \in \{\text{topics, style, etc.}\}$
 - speech recognition (signal processing) :
 $\phi_i \in \{\text{speaker gender, background noise, environment echo, etc.}\}$
- We will focus on vision problems
- We want to avoid hand-engineered solutions to these problems

Shallow Architecture Models

- A *shallow model* is a model with very few layers of computational units :



One hidden layer neural network

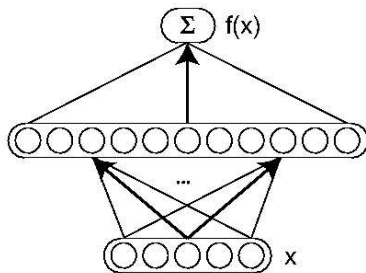


Kernel SVM

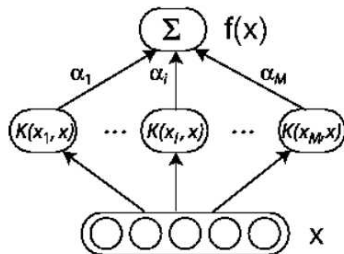
- To approximate a complex function, such models will need large (exponential) number of computational units

Shallow Architecture Models

- A *shallow model* is a model with very few layers of computational units :



One hidden layer neural network

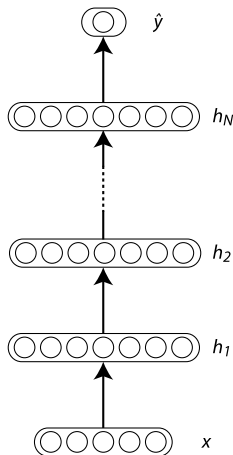


Kernel SVM

- To approximate a complex function, such models will need large (exponential) number of computational units

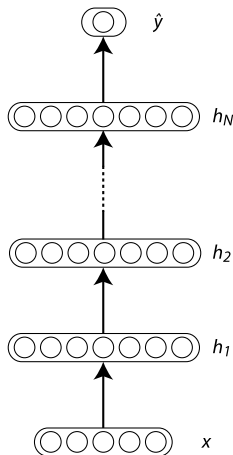
Deep Architecture Models

- A *deep architecture model* is such that its output is the result of the composition of many computational units
- Many layers potentially yield highly complex functions with a limited number of parameters
- The d dimensional parity function modeled with
 - $O(d2^d)$ parameters with Gaussian SVM
 - $O(d^2)$ parameters with a $O(\log_2 d)$ hidden layer neural network



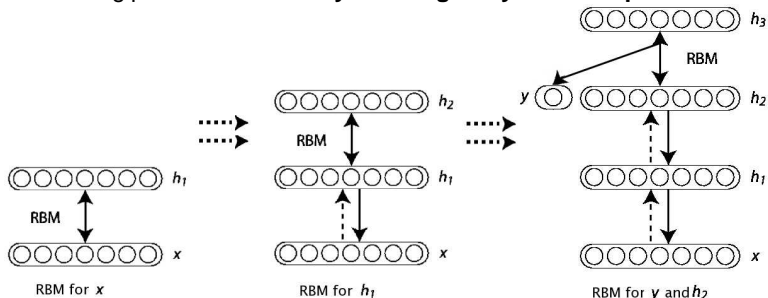
Deep Architecture Models

- A *deep architecture model* is such that its output is the result of the composition of many computational units
- Many layers potentially yield highly complex functions with a limited number of parameters
- The d dimensional parity function modeled with
 - $O(d2^d)$ parameters with Gaussian SVM
 - $O(d^2)$ parameters with a $O(\log_2 d)$ hidden layer neural network



Learning Deep Architecture Models (DBN)

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

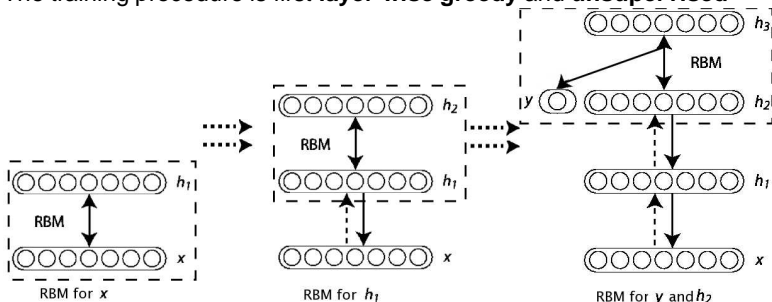


- Then the output of the model is **fine-tuned** on the supervised data

$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n \log \hat{p}(y_i | x_i, \theta)$$

Learning Deep Architecture Models (DBN)

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



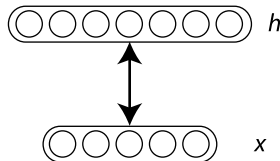
- Then the output of the model is **fine-tuned** on the supervised data

$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n \log \hat{p}(y_i | x_i, \theta)$$

Greedy Module : Restricted Boltzmann Machines

- RBM is an energy-based generative model :

$$p(X = x, H = h) \propto e^{-\mathcal{E}(x, h)} = e^{x'b + h'c + h'Wx}$$



- Inference is easy** : $p(H|X)$ factorizes ($H_i \perp H_j | X, \quad i \neq j$) :

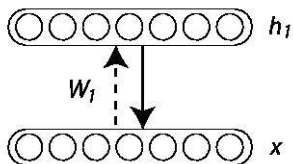
$$p(X^j = 1 | H = h) = \text{sigm}(b^j + \sum_i h^i W^{ij})$$

$$p(H^i = 1 | X = x) = \text{sigm}(c^i + \sum_j W^{ij} x^j).$$

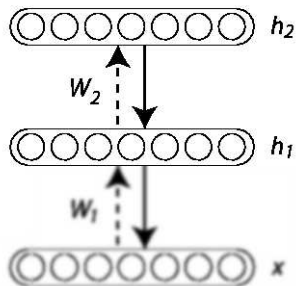
\Rightarrow no *explaining away*

- Training is easy** : Contrastive Divergence (Hinton 2002)

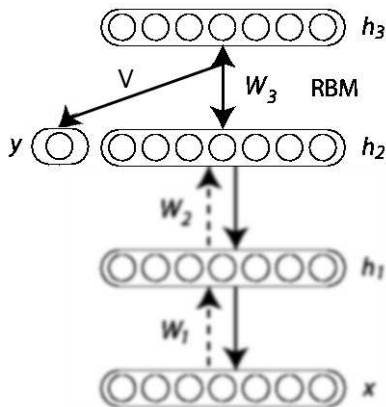
$$\frac{\partial \log p(x)}{\partial \theta} = - \sum_{h'} p(h'|x) \frac{\partial \mathcal{E}(x, h)}{\partial \theta} + \sum_{x', h'} p(x', h') \frac{\partial \mathcal{E}(x', h')}{\partial \theta}$$

Computing $\hat{p}(y|x)$ 

$$\hat{h}_1 = p(h_1|x) = \text{sigm}(c_1 + W_1x)$$

Computing $\hat{p}(y|x)$ 

$$\hat{h}_2 = p(h_2|\hat{h}_1) = \text{sigm}(c_2 + W_2\hat{h}_1)$$

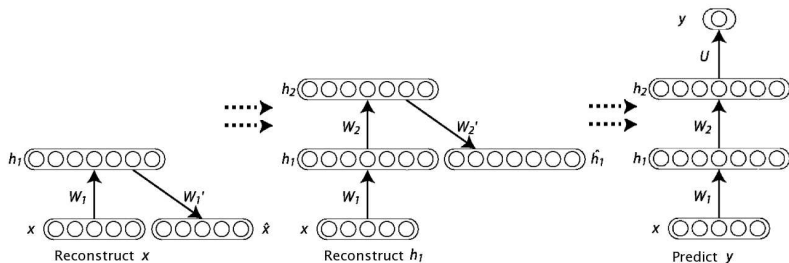
Computing $\hat{p}(y|x)$ 

$$\hat{p}(y|x) = \sum_{h_3} p(y, h^3 | \hat{h}_2)$$

$$\propto e^{c'_y y} \prod_k \left(1 + e^{V^k y + W_3^k \hat{h}_2 + c_3^k} \right)$$

Learning Deep Architecture Models (SAA)

- Instead of stacking RBMs, we can have Stacked Autoassociators (SAA)
- The training procedure is first **layer-wise greedy** and **unsupervised**

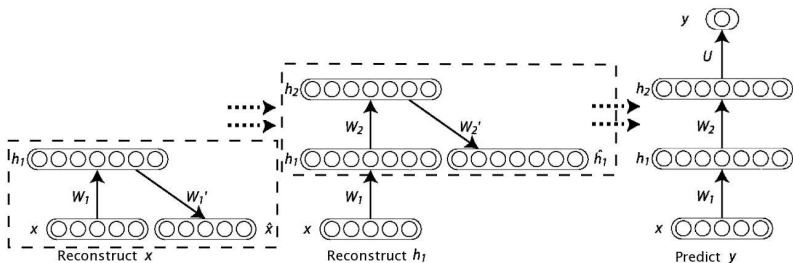


- Then the output of the model is **fine-tuned** on the supervised data

$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n \log \hat{p}(y_i | x_i, \theta)$$

Learning Deep Architecture Models (SAA)

- Instead of stacking RBMs, we can have Stacked Autoassociators (SAA)
- The training procedure is first **layer-wise greedy** and **unsupervised**



- Then the output of the model is **fine-tuned** on the supervised data

$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n \log \hat{p}(y_i | x_i, \theta)$$

Greedy Module : Autoassociators

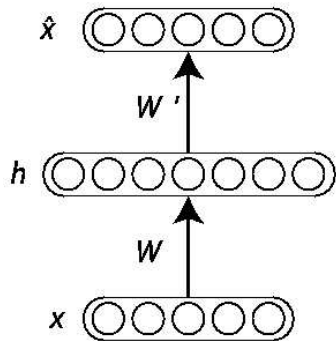
- Autoassociator is a neural network trained to reconstruct its input

$$\hat{x} = \text{sigm}(c + W' \text{sigm}(b + Wx))$$

- The reconstruction error is

$$R(x, \hat{x}) = - \sum_i x^i \log \hat{x}^i + (1 - x^i) \log(1 - \hat{x}^i)$$

- The neural network is trained using a gradient descent algorithm



Experimental setup

- We report results for the following models :
 - Support Vector Machine classifiers with polynomial (**SVM**_{poly}) and Gaussian (**SVM**_{rbf}) kernels
 - One hidden layer neural network (NNet)
 - Deep Belief Network (DBN-1 and DBN-3) and Stacked Autoassociator (SAA-3) with one or three hidden layers
- The validation set was used to do model selection and early stopping
- **SVM**_{poly} and **SVM**_{rbf} were retrained on the union of the training and validation set

Dataset Characteristics Summary

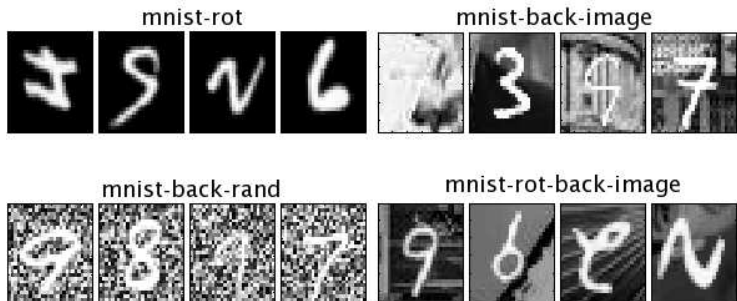
- Classification datasets on 28×28 pixel images
- All datasets have a training, validation and test split
- Training set size varies from 1000 to 10000 samples
- Validation set size varies from 200 to 2000 samples
- All datasets have a test set of size 50000.

Variations on Digit Recognition

- MNIST dataset with additional factors of variation
 - 1 Pick sample $(x, y) \in \mathcal{X}$ from the digit recognition dataset ;
 - 2 Create a perturbed version x^* of x according to some factors of variation ;
 - 3 Add (x^*, y) to a new dataset \mathcal{X}^* ;
 - 4 Go back to 1 until enough samples are generated.
- We generated the following datasets :

Dataset	Additional factors of variation
<i>mnist-rot</i>	rotation angle between 0 and 2π radians
<i>mnist-back-rand</i>	random background pixels between 0 and 255
<i>mnist-back-image</i>	random patch from 20 black and white images
<i>mnist-rot-back-image</i>	factors of <i>mnist-rot</i> and <i>mnist-back-image</i>

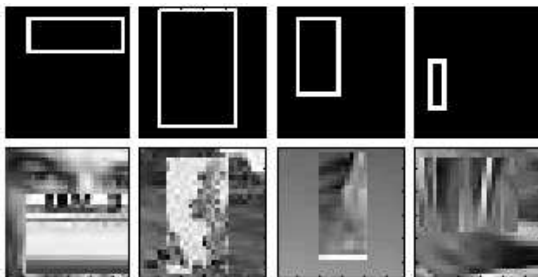
Variations on Digit Recognition (samples)



Dataset	SVM _{rbf}	SVM _{poly}	NNet	DBN-1	SAA-3	DBN-3
<i>mnist-basic</i>	3.03	3.69	4.69	3.94	3.46	3.11
<i>mnist-rot</i>	10.38	13.61	17.62	12.11	11.43	12.30
<i>mnist-back-rand</i>	14.58	16.62	20.04	9.80	11.28	6.73
<i>mnist-back-image</i>	22.61	24.01	27.41	16.15	23.00	16.31
<i>mnist-rot-back-image</i>	32.62	37.59	42.17	31.84	24.09	28.51

Discrimination between Tall and Wide Rectangles

- *rectangles*: the pixels corresponding to the border of the rectangle has a value of 255, 0 otherwise



- *rectangles-image*: the border and inside of the rectangles correspond to an image patch. A background patch is also sampled

Dataset	SVM _{rbf}	SVM _{poly}	NNet	DBN-1	SAA-3	DBN-3
<i>rectangles</i>	2.15	2.15	7.16	4.71	2.41	2.60
<i>rectangles-image</i>	24.04	24.05	33.20	23.69	24.05	22.50

Recognition of Convex Sets

- convex* contains images corresponding to convex and non convex sets of pixels



- The convex sets are intersections of random half-planes
- The non convex sets correspond to the union of a random number of convex sets, failing a convexity test

Dataset	SVM _{rbf}	SVM _{poly}	NNet	DBN-1	SAA-3	DBN-3
<i>convex</i>	19.13	19.82	32.25	19.92	18.41	18.63

Results Summary

Dataset	SVM _{rbf}	SVM _{poly}	NNet	DBN-1	SAA-3	DBN-3
<i>mnist-basic</i>	3.03	3.69	4.69	3.94	3.46	3.11
<i>mnist-rot</i>	10.38	13.61	17.62	12.11	11.43	12.30
<i>mnist-back-rand</i>	14.58	16.62	20.04	9.80	11.28	6.73
<i>mnist-back-image</i>	22.61	24.01	27.41	16.15	23.00	16.31
<i>mnist-rot-back-image</i>	32.62	37.59	42.17	31.84	24.09	28.51
<i>rectangles</i>	2.15	2.15	7.16	4.71	2.41	2.60
<i>rectangles-image</i>	24.04	24.05	33.20	23.69	24.05	22.50
<i>convex</i>	19.13	19.82	32.25	19.92	18.41	18.63

TABLE: Results on the benchmark for problems with factors of variation (in percentages). The best performance as well as those with overlapping confidence intervals are marked in bold.

Investigation of the “background effect”

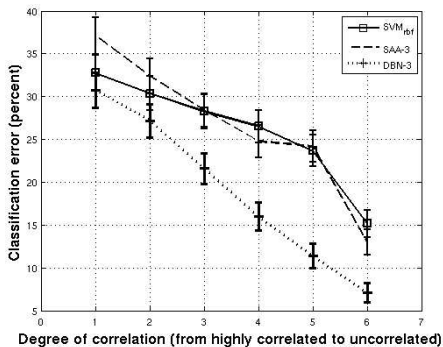
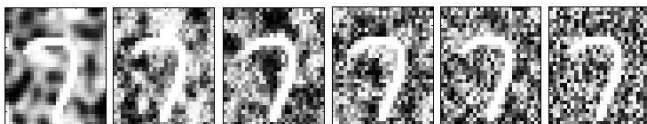


FIG.: Classification error of SVM_{rbf} , SAA-3 and DBN-3 on MNIST examples with progressively less pixel correlation in the background.

Conclusion

- In general, models with deep architectures either perform as well or outperform other models
- There are still challenges in scaling the current algorithms to problems with very complex input distribution
- Datasets and experimental details can be found on our public wiki page :

`http://www.iro.umontreal.ca/~lisa/ptwiki/`

Future Work

- Address the “focus problem” of greedy layer-wise unsupervised training
- Develop learning algorithms that make better use of the capacity of the model or models appropriate for more complex input distributions
- Develop models and algorithms with less hyper-parameters

THANK YOU!

Hinton, G. E., Osindero, S., & Teh, Y. (2006). A fast learning algorithm for deep belief nets. *Neural Computation*, 18, 1527–1554.