

# A Review of Deep Models for Density Estimation

Andrew McHutchon & Mark van der Wilk

Thursday 13 March, 2014

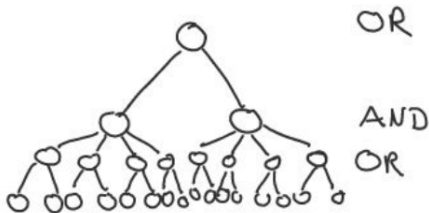
- ▶ Justifying deep architectures (intuition and a theorem)
- ▶ Where it started: Deep RBMs
- ▶ The Bayesian approach: Deep GPs
- ▶ A different approach: Deep Density Models
- ▶ A hybrid: Deep Latent Gaussian Models
- ▶ Tying it together

- ▶ Some functions cannot be “efficiently” represented by shallow architectures
- ▶ Shallow architectures generalise locally, while deep representations allow more global generalisation
- ▶ Multi-task learning. Low level features can be useful for several tasks.
- ▶ Empirical performance

[Bengio, 2009]

## Theorem

A monotone weighted threshold circuit of depth  $k - 1$  computing a function  $f_k \in \mathcal{F}_{k,N}$  has at least size  $2^{cN}$  for some constant  $c > 0$  and  $N > N_0$ . A depth  $k$  circuit needs only polynomial size.  
[Håstad and Goldman, 1991]



- ▶ Proof is limited to one type of NN and one class of functions... But it does hint that there can be some advantage to using deep architectures.
- ▶ Is this relevant to “infinite capacity” non-parametric models (GP’s)?

Bengio: *Shallow architectures (kernel learning) generalise locally, while deep representations allow more global generalisation.*

To a certain extent, this is a simplification of something we already know.

- ▶ Bengio's point: Squared Exponential is very widespread in use, but provides only local generalisation.
- ▶ Good features can greatly improve generalisation.
- ▶ A feature mapping can be incorporated into a the kernel (like the periodic kernel).
- ▶ Finding good features/kernels is hard
- ▶ Learning deep representations may be a way of doing this.

Bengio: *Shallow architectures (kernel learning) generalise locally, while deep representations allow more global generalisation.*

To a certain extent, this is a simplification of something we already know.

- ▶ Bengio's point: Squared Exponential is very widespread in use, but provides only local generalisation.
- ▶ Good features can greatly improve generalisation.
- ▶ A feature mapping can be incorporated into a the kernel (like the periodic kernel).
- ▶ Finding good features/kernels is hard
- ▶ Learning deep representations may be a way of doing this.

Bengio: *Shallow architectures (kernel learning) generalise locally, while deep representations allow more global generalisation.*

To a certain extent, this is a simplification of something we already know.

- ▶ Bengio's point: Squared Exponential is very widespread in use, but provides only local generalisation.
- ▶ Good features can greatly improve generalisation.
- ▶ A feature mapping can be incorporated into a the kernel (like the periodic kernel).
- ▶ Finding good features/kernels is hard
- ▶ Learning deep representations may be a way of doing this.

Bengio: *Shallow architectures (kernel learning) generalise locally, while deep representations allow more global generalisation.*

To a certain extent, this is a simplification of something we already know.

- ▶ Bengio's point: Squared Exponential is very widespread in use, but provides only local generalisation.
- ▶ Good features can greatly improve generalisation.
- ▶ A feature mapping can be incorporated into a the kernel (like the periodic kernel).
- ▶ Finding good features/kernels is hard
- ▶ Learning deep representations may be a way of doing this.



Bengio: *Shallow architectures (kernel learning) generalise locally, while deep representations allow more global generalisation.*

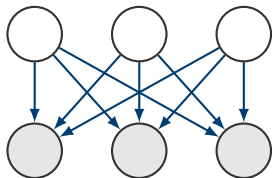
To a certain extent, this is a simplification of something we already know.

- ▶ Bengio's point: Squared Exponential is very widespread in use, but provides only local generalisation.
- ▶ Good features can greatly improve generalisation.
- ▶ A feature mapping can be incorporated into a the kernel (like the periodic kernel).
- ▶ Finding good features/kernels is hard
- ▶ Learning deep representations may be a way of doing this.

Bengio: *Shallow architectures (kernel learning) generalise locally, while deep representations allow more global generalisation.*

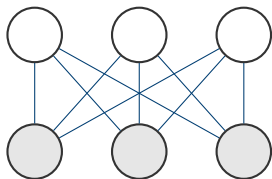
To a certain extent, this is a simplification of something we already know.

- ▶ Bengio's point: Squared Exponential is very widespread in use, but provides only local generalisation.
- ▶ Good features can greatly improve generalisation.
- ▶ A feature mapping can be incorporated into a the kernel (like the periodic kernel).
- ▶ Finding good features/kernels is hard
- ▶ Learning deep representations may be a way of doing this.



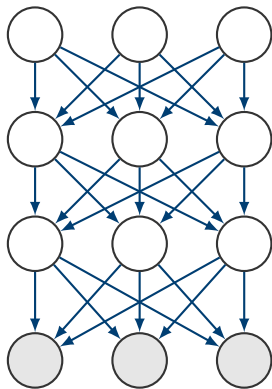
## Belief net:

- ▶ Directed acyclic graph
- ▶ Sigmoid belief net [Neal (1992)]
  - ▶ Binary nodes
  - ▶  $p(\mathbf{y}, \mathbf{h}) = \prod_{i=1}^D \text{Ber}(y_i | \sigma(\mathbf{w}^T \mathbf{h})) \prod_{j=1}^E \text{Ber}(h_j | w_{2j})$



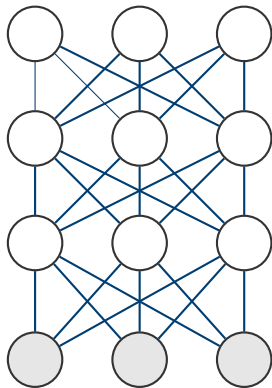
## Restricted Boltzmann machine:

- ▶ Undirected bipartite graph
- ▶ Binary nodes
- ▶  $p(\mathbf{y}, \mathbf{h}) = \frac{\exp(\mathbf{h}^T \mathbf{W} \mathbf{y} + \mathbf{a}^T \mathbf{y} + \mathbf{b}^T \mathbf{h})}{Z}$

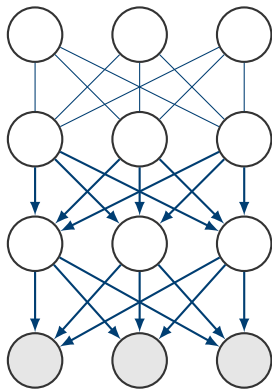


**Deep directed net:** Training very difficult due to “explaining away” effect

- ▶ Posterior on hidden nodes does not factorise
- ▶ MCMC - slow
- ▶ Variational - too approximate

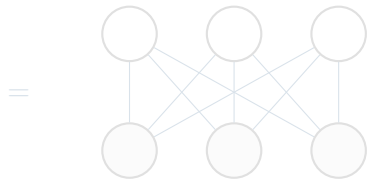
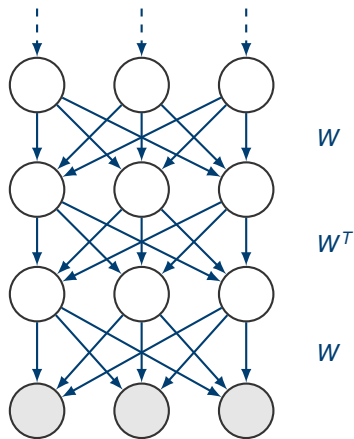


**Deep Boltzmann machine:**  
Sulakhutinov & Hinton, AISTATS  
(2009)

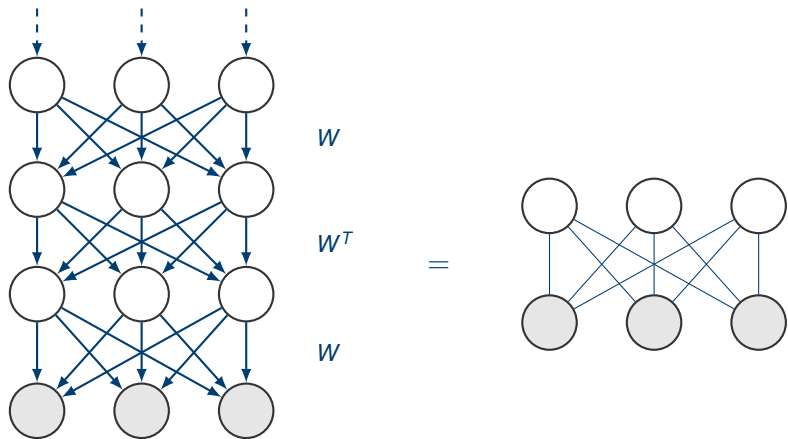


## Deep Belief net:

- ▶ Hinton, Osindero & Teh, Neural Computation (2006)
- ▶ RBM acts as a “complementary prior” to remove “explaining away”



Posterior on hidden nodes  
factorises



Posterior on hidden nodes  
factorises

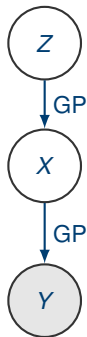


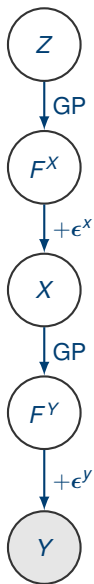
## Training:

1. Iteratively train each layer in a greedy manner
2. Fine tuning with and up-down algorithm on subset of data
3. Cross validation step
4. Further training

- ▶ Gaussian Process Latent Variable Model  
[Lawrence (2004)]
- ▶ Bayesian Gaussian Process Latent Variable Models  
[Titsias and Lawrence (2010)]
- ▶ Deep Gaussian Processes  
[Damianou and Lawrence (2013)]







$$F^X \sim \mathcal{GP}(\mathbf{0}, k^X(Z, Z)) \quad (1)$$

$$\mathbf{x}_n = \mathbf{f}_n^X + \epsilon_n^X \quad (2)$$

$$F^Y \sim \mathcal{GP}(\mathbf{0}, k^X(X, X)) \quad (3)$$

$$\mathbf{y}_n = \mathbf{f}_n^Y + \epsilon_n^Y \quad (4)$$

$$\log p(Y) = \log \int p(Y | F^Y) p(F^Y | X) p(X | F^X) p(F^X | Z) p(Z) \quad (5)$$

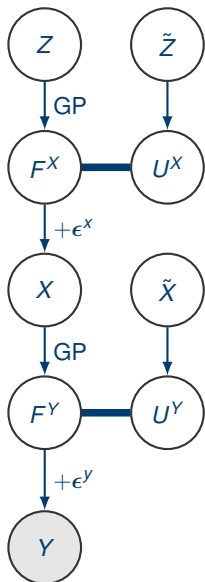
$$\mathcal{F} = \int Q \log \frac{p(Y | F^Y) p(F^Y | X) p(X | F^X) p(F^X | Z) p(Z)}{Q} \quad (6)$$

$$\mathcal{F} \leq \log p(Y) \quad (7)$$

$$\log p(Y) = \log \int p(Y | F^Y) p(F^Y | X) p(X | F^X) p(F^X | Z) p(Z) \quad (5)$$

$$\mathcal{F} = \int Q \log \frac{p(Y | F^Y) p(F^Y | X) p(X | F^X) p(F^X | Z) p(Z)}{Q} \quad (6)$$

$$\mathcal{F} \leq \log p(Y) \quad (7)$$



$$p(F^X | Z) = p(F^X | U^X, Z, \tilde{Z}) p(U^X | \tilde{Z}) \quad (8)$$

$$p(F^Y | X) = p(F^Y | U^Y, X, \tilde{X}) p(U^Y | \tilde{X}) \quad (9)$$



$$Q = p(F^Y | U^Y, X)q(U^Y)q(X)p(F^Z | U^Z, Z)q(U^Z)q(Z) \quad (10)$$

$$\mathcal{F} = \int Q \log \frac{p(Y | F^Y)p(F^Y | X)p(X | F^X)p(F^X | Z)p(Z)}{p(F^Y | U^Y, X)q(U^Y)q(X)p(F^Z | U^Z, Z)q(U^Z)q(Z)}$$

(11)

$$Q = p(F^Y | U^Y, X)q(U^Y)q(X)p(F^Z | U^Z, Z)q(U^Z)q(Z) \quad (10)$$

$$\mathcal{F} = \int Q \log \frac{p(Y | F^Y)p(F^Y | X)p(X | F^X)p(F^X | Z)p(Z)}{p(F^Y | U^Y, X)q(U^Y)q(X)p(F^Z | U^Z, Z)q(U^Z)q(Z)}$$

(11)

$$Q = p(F^Y | U^Y, X)q(U^Y)q(X)p(F^Z | U^Z, Z)q(U^Z)q(Z) \quad (10)$$

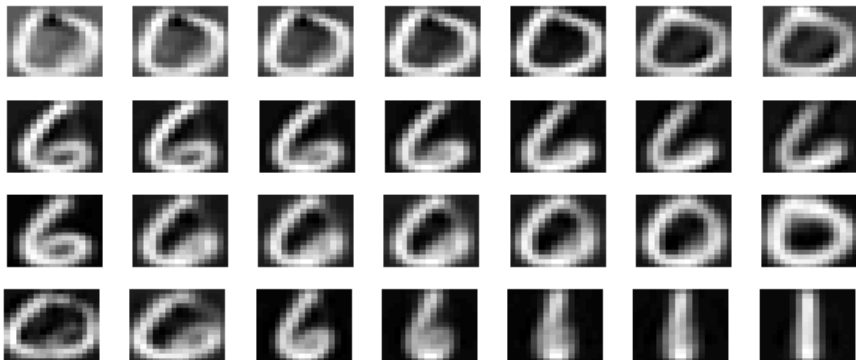
$$\mathcal{F} = \int Q \log \frac{p(Y | F^Y)p(F^Y | X)p(X | F^X)p(F^X | Z)p(Z)}{p(F^Y | U^Y, X)q(U^Y)q(X)p(F^Z | U^Z, Z)q(U^Z)q(Z)}$$

$$= \int Q \log \frac{p(Y | F^Y)p(U^Y)p(X | F^X)p(U^X)p(Z)}{q(U^Y)q(X)q(U^Z)q(Z)}$$

(11)

$$Q = p(F^Y | U^Y, X)q(U^Y)q(X)p(F^Z | U^Z, Z)q(U^Z)q(Z) \quad (10)$$

$$\begin{aligned} \mathcal{F} &= \int Q \log \frac{p(Y | F^Y)p(F^Y | X)p(X | F^X)p(F^X | Z)p(Z)}{p(F^Y | U^Y, X)q(U^Y)q(X)p(F^Z | U^Z, Z)q(U^Z)q(Z)} \\ &= \int Q \log \frac{p(Y | F^Y)p(U^Y)p(X | F^X)p(U^X)p(Z)}{q(U^Y)q(X)q(U^Z)q(Z)} \\ &= g_Y + r_X + H(q(X)) - KL(q(Z) \parallel p(Z)) \end{aligned} \quad (11)$$



- ▶ USPS digits {0, 1, 6}
- ▶ 5 layers
- ▶ 150 training data points!!!

## Parameters

- ▶  $M$   $D$ -dimensional pseudo points per layer
- ▶  $D + 2$  hyperparameters per layer
- ▶  $D + 1$  variational parameters per layer

## Complexity

- ▶  $\mathcal{O}(NM^2L)$  training
- ▶  $\mathcal{O}(M^2L)$  generative (with precomputations)

Introduced by Rippel and Adams [2013].

Main idea:

- ▶ Find an *invertible* mapping  $\mathbf{y} = \mathbf{f}(\mathbf{x})$  such that  $p(\mathbf{x})$  has independent dimensions with *given marginals*:

$$p(\mathbf{x}) = \prod_{k=1}^K p_{X_k}(x_k)$$

- ▶ i.e. transform the data such that the distribution in the “representation space” is known.

Consequences:

- ▶ This makes the density over  $\mathbf{y}$  easy to calculate:

$$p(\mathbf{y}) = \prod_{k=1}^K p_{X_k}(\mathbf{f}^{-1}(\mathbf{y})) \left| \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{f}^{-1}(\mathbf{y})}$$

- ▶ Normalised and fast to compute  $p(\mathbf{y}^*)!$

Introduced by Rippel and Adams [2013].

Main idea:

- ▶ Find an *invertible* mapping  $\mathbf{y} = \mathbf{f}(\mathbf{x})$  such that  $p(\mathbf{x})$  has independent dimensions with *given marginals*:

$$p(\mathbf{x}) = \prod_{k=1}^K p_{X_k}(x_k)$$

- ▶ i.e. transform the data such that the distribution in the “representation space” is known.

Consequences:

- ▶ This makes the density over  $\mathbf{y}$  easy to calculate:

$$p(\mathbf{y}) = \prod_{k=1}^K p_{X_k}(\mathbf{f}^{-1}(\mathbf{y})) \left| \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{f}^{-1}(\mathbf{y})}$$

- ▶ Normalised and fast to compute  $p(\mathbf{y}^*)!$



The decoder  $\mathbf{f}(\mathbf{x})$  and its inverse  $\mathbf{g}(\mathbf{y}) \approx \mathbf{f}^{-1}(\mathbf{y}) = \mathbf{x}$  are both parameterised by deep neural networks:

$$\mathbf{f}_{\Theta}(\mathbf{x}) = \sigma(\Omega_M \cdot + \omega_M) \circ \sigma(\Omega_{M-1} \cdot + \omega_{M-1}) \circ \cdots \circ \sigma(\Omega_1 \mathbf{x} + \omega_1)$$

$$\mathbf{g}_{\Psi}(\mathbf{y}) = \sigma(\Gamma_M \cdot + \gamma_M) \circ \sigma(\Gamma_{M-1} \cdot + \gamma_{M-1}) \circ \cdots \circ \sigma(\Gamma_1 \mathbf{y} + \gamma_1)$$

We find  $\mathbf{g} \approx \mathbf{f}^{-1}$  separately because this has regularisation benefits (discussed shortly).

The cost function:

$$C(\Theta, \Psi) = \mu_{\mathcal{D}} \mathcal{D}(\Psi) + \mu_{\mathcal{I}} \mathcal{I}(\Theta) + \mu_{\mathcal{R}} \mathcal{R}(\Theta, \Psi)$$

- ▶  $\mathcal{D}(\Psi)$ : Divergence penalty
- ▶  $\mathcal{I}(\Theta)$ : Invertability penalty
- ▶  $\mathcal{R}(\Theta, \Psi)$ : Reconstruction loss

Forces the model to distribute mass of the data in the representation space similarly to the desired distribution  $p(\mathbf{x})$ .

$$\mathcal{D}(\Psi) = \left( \frac{1}{K} \sum_{k=1}^K T(\hat{p}_{X_k}(\cdot) \| p_{X_k}(\cdot)) \right) + \left( \frac{1}{N} \sum_{n=1}^N T(\hat{p}_X(\mathbf{x}_n) \| p_X(\mathbf{x}_n)) \right)$$

$$\hat{p}_{X_k}(x) = \frac{1}{N} \sum_{n=1}^N \delta(x - [\mathbf{x}_n]_k)$$

$$\mathcal{I}(\Theta) = \frac{1}{M} \sum_{m=1}^M \log \left( \frac{\lambda_{\max}(\Omega_m)}{\lambda_{\min}(\Omega_m)} \right)$$

Ensures invertibility, so  $p(\mathbf{y})$  is well defined.

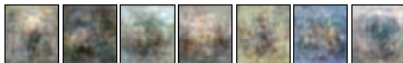
- ▶ Ensures that  $\mathbf{g}_\psi(\mathbf{y}) \approx \mathbf{f}_\theta^{-1}(\mathbf{y})$  - Information is preserved.
- ▶ As a consequence, marginals are forced to become independent.

$$\begin{aligned} \mathcal{H}(p_{\mathbf{Y}}(\cdot)) &= \sum_{k=1}^K \mathcal{H}(p_{X_k}(\cdot)) + \\ &\quad - D\left(\prod_{k=1}^K p_{X_k}(\cdot) \parallel p_{\mathbf{X}}(\cdot)\right) + \\ &\quad + \mathbb{E} \left[ \log \left| \frac{\partial f_{\theta}(\cdot)}{\partial \mathbf{y}} \right| \right] \end{aligned}$$

- ▶ 1.6% error rate on MNIST



(a) Generations from a 3-layer DDM on MNIST, with diversified marginals Beta ( $\cdot$ ; 0.015, 0.8).



(b) Generations from a 3-layer DDM on CIFAR-10, with diversified marginals Beta ( $\cdot$ ; 0.5, 3).

Figure 4: Generations from DDMs.

Claims to be:

- ▶ Deep - for capturing higher moments of the data
- ▶ Non-linear - Allowing for complex structure in the data
- ▶ Fast - in terms of sampling fantasy data
- ▶ Scalable - in terms of the size of training set

Also provides a nice link to autoencoders.

## Generative process

$$\xi_l \sim \mathcal{N}(\xi_l | 0, I) \quad (12)$$

$$\mathbf{h}_L = G_L \xi_L \quad (13)$$

$$\mathbf{h}_l = T_l(\mathbf{h}_{l+1}) + G_l \xi_l \quad (14)$$

$$\mathbf{v} \sim p(\mathbf{v} | T_0(\mathbf{h}_1)) \quad (15)$$

- ▶ Transformation  $T_l$  is a multilayer perceptron
- ▶ Parameters  $\theta$  contains the  $\{G_l\}$  and the MLP parameters



- ▶ Variationally integrate out  $\xi$ :  
$$\mathcal{L}(V) = -D_{KL}(q(\xi)||p(\xi)) + \mathbb{E}_q [\log p(V|\xi, \theta^g)]$$
- ▶ Need gradients w.r.t.  $\theta^g$  and the variational parameters

The trick for the variational parameters:

- ▶  $\nabla_{\mu_j} \mathbb{E}_q [\log p(V|\xi, \theta^g)] = \mathbb{E} [\nabla_{\xi_j} \log p(\mathbf{v}|\mathbf{h}(\xi))]$
- ▶ Backpropagation gives the gradients given  $\xi$
- ▶ Sampling can give an unbiased estimate of the gradient

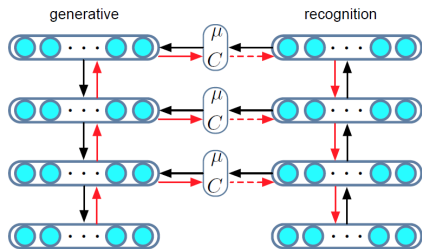
- ▶ Variationally integrate out  $\xi$ :  
$$\mathcal{L}(V) = -D_{KL}(q(\xi)||p(\xi)) + \mathbb{E}_q [\log p(V|\xi, \theta^g)]$$
- ▶ Need gradients w.r.t.  $\theta^g$  and the variational parameters

The trick for the variational parameters:

- ▶  $\nabla_{\mu_j} \mathbb{E}_q [\log p(V|\xi, \theta^g)] = \mathbb{E} [\nabla_{\xi_j} \log p(\mathbf{v}|\mathbf{h}(\xi))]$
- ▶ Backpropagation gives the gradients given  $\xi$
- ▶ Sampling can give an unbiased estimate of the gradient

## Training buzzwords: “Learning to learn Variational Stochastic Back-Propagation”

- ▶ Estimate  $\mu_l$  &  $C_l$  using the recognition model
- ▶ Sample  $\xi$  for each layer so the  $\mathbf{h}$ 's can be calculated
- ▶ Calculate the gradient estimates w.r.t.  $\theta^g$
- ▶ Calculate the gradient estimates w.r.t.  $\theta^r$



Autoencoder objective function:

$$\mathcal{L}(\theta) = \frac{1}{n} \sum_{V \sim p(V), \tilde{V} \sim \mathcal{L}(\tilde{V}|V)} \lambda_n \Omega(\theta, V, \tilde{V}) - \log P_\theta(V|\tilde{V})$$

Deep Latent Gaussian Model:

$$\mathcal{L}(V) = -D_{KL}(q(\xi)||p(\xi)) + \mathbb{E}_q [\log p(V|\xi, \theta^g)]$$

- ▶ Autoencoders aim to maximise the probability of reconstructing a corrupted sample
- ▶ The noisy sample  $\tilde{V}$  is analogous to the latent Gaussians  $\xi$
- ▶ The learned posterior  $q(\xi|\mathbf{v})$  is analogous to the encoder

- ▶ All: Attempt to map the observed space to some latent space with an easy density
  - ▶ Some proofs about why this may be good: Gaussianisation
- ▶ Latent dimension different for all
- ▶ Inference varies wildly



Y. Bengio. Learning deep architectures for ai. *Found. Trends Mach. Learn.*, 2(1):1–127, January 2009. ISSN 1935-8237. doi: 10.1561/22000000006. URL

<http://dx.doi.org/10.1561/22000000006>.

Johan Håstad and Mikael Goldmann. On the power of small-depth threshold circuits. *computational complexity*, 1(2):113–129, 1991. ISSN 1016-3328. doi: 10.1007/BF01272517. URL

<http://dx.doi.org/10.1007/BF01272517>.

Oren Rippel and Ryan Prescott Adams. High-dimensional probability estimation with deep density models, 2013.