



UNIVERSIDAD
COMPLUTENSE
MADRID

université
PARIS-SACLAY

Data modeling with Energy Based Models

Beatriz Seoane

LISN Paris-Saclay University

Acknowledgments

Aurélien Decelle
Giovanni Catania
Alfonso Navas
Lorenzo Rosset



Nicolas Béreux
Cyril Furtlehner





Plan for the lecturers

- Class 1: **Introduction** to Energy Based Models
- Class 2: **Interpretability**. How can we learn from trained networks?
- Class 3: **Training optimization, the role of MCMC**. How can we improve the training mechanisms by understanding their physics?



Plan for the lecturers

- Class 1: **Introduction** to Energy Based Models
 - Generative approach
 - Introduction to Energy-Based Models
 - The Restricted Boltzmann Machine (RBM)
 - Maximum likelihood training
 - Generation
 - Why I think RBMs are a cool tool

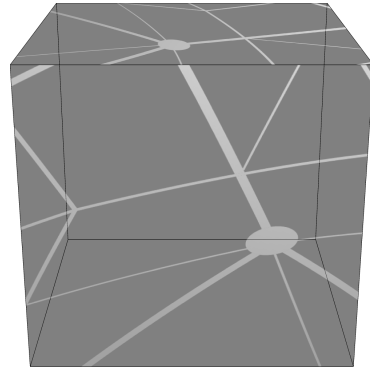


General definitions

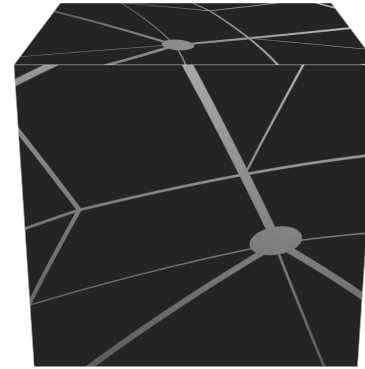
Introduction : Generative approach

0
1
2
3
4
5
6
7
8
9

training



generating

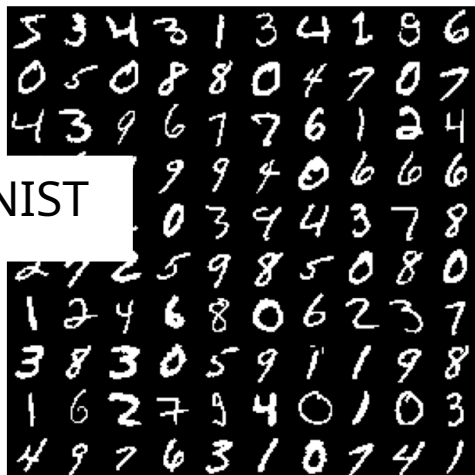


1
6
9
6
4
7
9
8
7
5

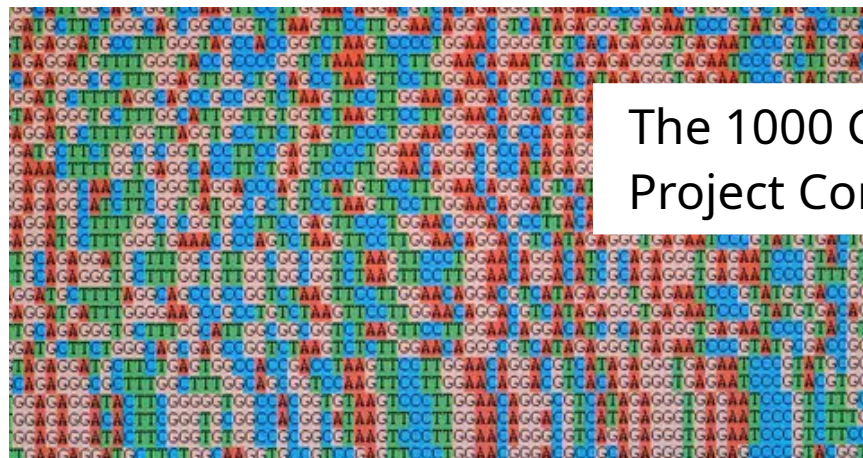
- **Energy based models (RBMs, Generative Convnets)**
- **Diffusion models**, normalizing flows, score based
- Variational AutoEncoder (VAE)
- Generative Adverarial Network (GAN)
- Autoregressive methods

Introduction : generative approach

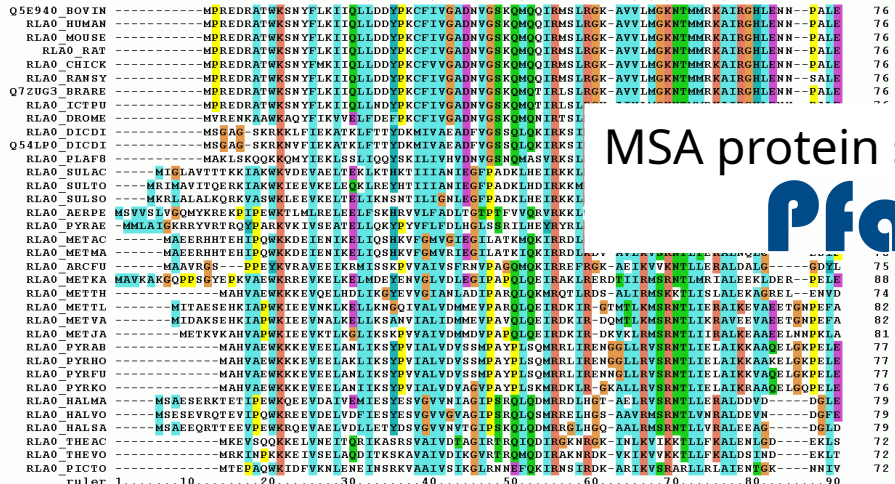
MNIST



CELEBA



The 1000 Genomes Project Consortium



MSA protein sequences



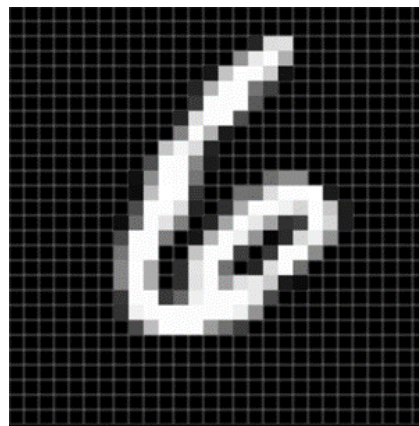
Data

$$\mathcal{D} = \left\{ \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)} \right\}$$

m -th entry $\mathbf{x}^{(m)} = \begin{bmatrix} x_1^{(m)} \\ \vdots \\ x_N^{(m)} \end{bmatrix}$

3	8	6	9	6	4	5	3	8	4	5	2	3	8	4	8
1	5	0	5	9	7	4	1	0	3	0	6	2	9	9	4
1	3	6	8	0	7	7	6	8	9	0	3	8	3	7	7
8	4	4	1	2	9	8	1	1	0	6	6	5	0	1	1

M : # of examples in the data set

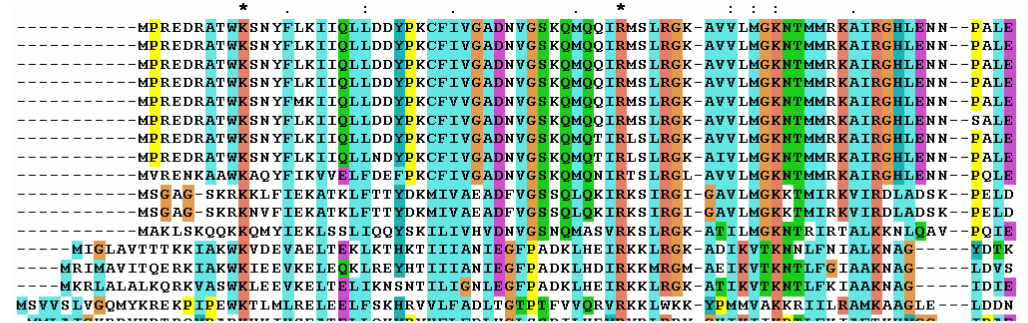


$N = 28 \times 28$
pixels

Data

$$\mathcal{D} = \left\{ \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)} \right\}$$

m -th entry $\mathbf{x}^{(m)} = \begin{bmatrix} x_1^{(m)} \\ \vdots \\ x_N^{(m)} \end{bmatrix}$



M: # of sequences in a protein family



$$N = L_{MSA} \quad \text{Amino-acids}$$

Goal: Create synthetic sequences

Data

$$\mathcal{D} = \left\{ \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)} \right\}$$



m-th entry $\mathbf{x}^{(m)} = \begin{bmatrix} x_1^{(m)} \\ \vdots \\ x_N^{(m)} \end{bmatrix}$

$\in \mathbb{R}^N$	continuous
$\in [0, 1]^N$	binary
$\in [G, A \dots, N, Q, -]^N$	categorical

Data distribution

$$\mathcal{D} = \left\{ \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)} \right\}$$

Underlying assumption

i. i. d. realizations of a random variable

$$\mathbf{X} \sim P_{\text{data}} \quad (\text{Generally unknown})$$

Empirical data distribution

$$\mathcal{D} = \left\{ \mathbf{x}_d^{(1)}, \dots, \mathbf{x}_d^{(M)} \right\}$$

Underlying assumption

i. i. d. realizations of a random variable

$$\mathbf{X}_d \sim P_{\text{data}} \quad (\text{Generally unknown})$$

Empirical distribution

$$p_{\mathcal{D}}(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M \delta \left(\mathbf{x} - \mathbf{x}_d^{(m)} \right) \xrightarrow{\text{Large } M} p_{\text{data}}(\mathbf{x})$$

Empirical data distribution

$$\mathcal{D} = \left\{ \mathbf{x}_d^{(1)}, \dots, \mathbf{x}_d^{(M)} \right\}$$

Underlying assumption

i. i. d. realizations of a random variable

$$\mathbf{X}_d \sim P_{\text{data}} \quad (\text{Generally unknown})$$

Empirical distribution

$$p_{\mathcal{D}}(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M \delta \left(\mathbf{x} - \mathbf{x}_d^{(m)} \right) \xrightarrow{\text{Large } M} p_{\text{data}}(\mathbf{x})$$

Generative models : $p_{\theta}(\mathbf{x})$ θ



Energy-based models

Energy based models (EBMs)

Hinton, Hopfield, LeCun, Bengio

Empirical

Model

$$p_{\mathcal{D}}(\mathbf{x}) \sim p_{\theta}(\mathbf{x}) = \frac{e^{-E_{\theta}(\mathbf{x})}}{Z_{\theta}}$$

Gibbs-Boltzmann distribution

$E_{\theta}(\mathbf{x})$ energy function

$$Z_{\theta} = \int d\mathbf{x} e^{-E_{\theta}(\mathbf{x})} \quad \text{Partition function}$$

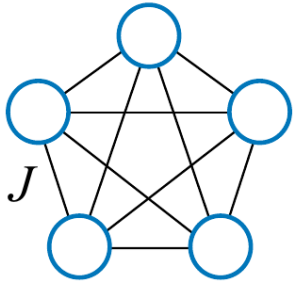
Learning : adjust the parameters θ so that the dataset configurations are **typical** configurations of the model.

Energy based models (EBMs)

Hinton, Hopfield, LeCun, Bengio

Boltzmann Machines (Ising/Hopfield/Potts models)

- Ackley, D. H., Hinton, G. E., & Sejnowski, T. J. (1985). *A learning algorithm for Boltzmann machines*. *Cognitive science*, 9(1), 147-169.



$$E_{J,h}(\mathbf{x}) = -\mathbf{x}^\top J \mathbf{x} - \mathbf{h}^\top \mathbf{x}$$

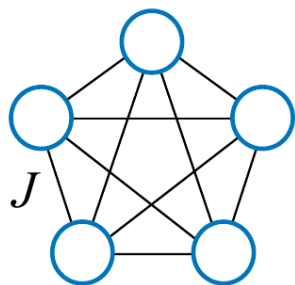
Pairwise interactions

$$E_{\theta}(\mathbf{x})$$

Energy based models (EBMs)

Hinton, Hopfield, LeCun, Bengio

- Ising/Hopfield/Potts models



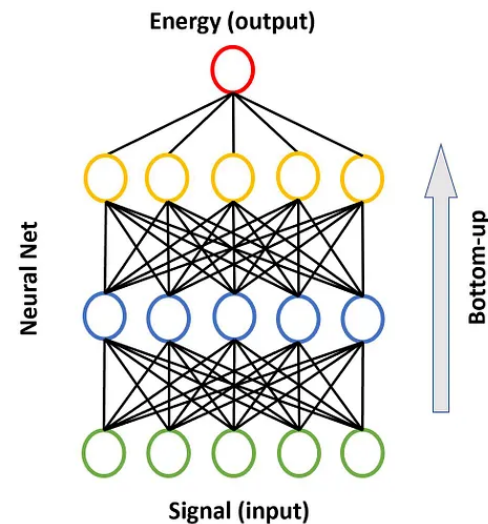
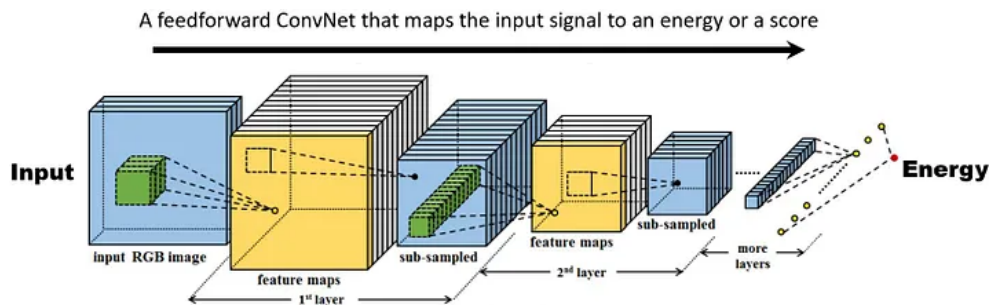
$$E_{J,h}(\mathbf{x}) = -\mathbf{x}^\top J \mathbf{x} - \mathbf{h}^\top \mathbf{x}$$

$$E_{\theta}(\mathbf{x})$$

- Generative ConvNets

- LeCun, Y., Chopra, S., Hadsell, R., Ranzato, M., & Huang, F. (2006). *A tutorial on energy-based learning*.

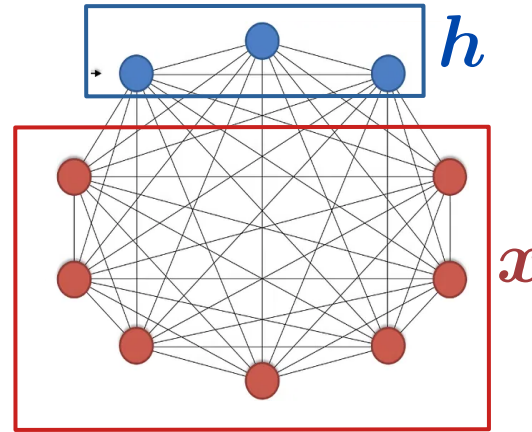
- Xie, J., Lu, Y., Zhu, S. C., & Wu, Y. (2016). *A theory of generative convnet*.



Models with hidden variables

- Boltzmann Machines

- Ackley, D. H., Hinton, G. E., & Sejnowski, T. J. (1985). *A learning algorithm for Boltzmann machines*. *Cognitive science*, 9(1), 147-169.



visible variables
 $\mathcal{E}(x, h; \theta)$

Latent variables
Encode correlations

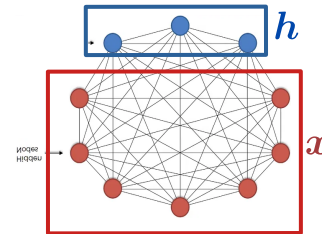
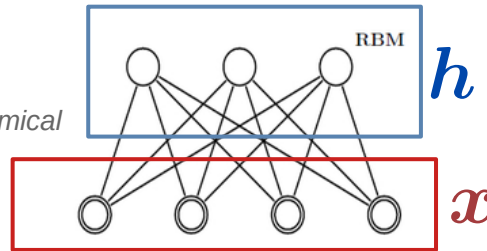
Models with hidden variables

- Boltzmann Machines

- Ackley, D. H., Hinton, G. E., & Sejnowski, T. J. (1985). *A learning algorithm for Boltzmann machines*. *Cognitive science*, 9(1), 147-169.

- Restricted Boltzmann Machine

- Smolensky, P. (1986). *Information processing in dynamical systems: Foundations of harmony theory*.



$$\mathcal{E}(\mathbf{x}, \mathbf{h}; \boldsymbol{\theta})$$

visible variables

Latent variables
Encode correlations

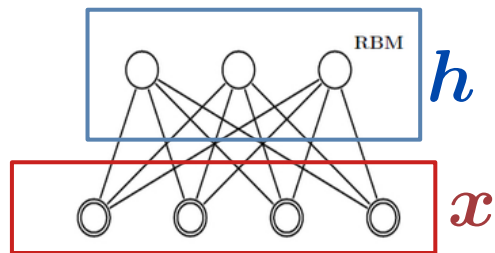
Hidden layer : interactions

Visible layer : data

Models with hidden variables

- Restricted Boltzmann Machine

- Smolensky, P. (1986). *Information processing in dynamical systems: Foundations of harmony theory.*



visible variables
 $\mathcal{E}(\mathbf{x}, \mathbf{h}; \theta)$
Latent variables
Encode correlations

$$\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h}) = -\mathbf{x}^{\top} W \mathbf{h} - \zeta^{\top} \mathbf{x} - \eta^{\top} \mathbf{h}$$

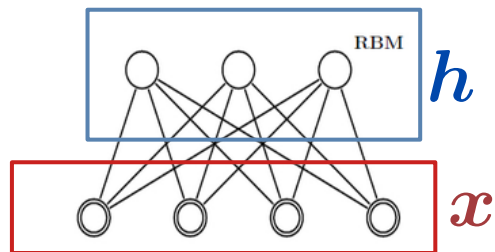
$$\theta = \{W, \zeta, \eta\}$$

$$p_{\theta}(\mathbf{x}) = \frac{e^{-E_{\theta}(\mathbf{x})}}{Z_{\theta}} = \frac{\sum_{\mathbf{h}} e^{-\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h})}}{Z_{\theta}} = \frac{e^{\sum_i x_i \zeta_i}}{Z_{\theta}} \prod_{a=1}^{N_h} \sum_{h_a=0}^1 e^{\sum_i x_i W_{ia} h_a + \eta_a h_a}$$

Models with hidden variables

- Restricted Boltzmann Machine

- Smolensky, P. (1986). *Information processing in dynamical systems: Foundations of harmony theory.*



visible variables
 $\mathcal{E}(\mathbf{x}, \mathbf{h}; \boldsymbol{\theta})$
Latent variables
Encode correlations

$$\mathcal{E}_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{h}) = -\mathbf{x}^{\top} W \mathbf{h} - \boldsymbol{\zeta}^{\top} \mathbf{x} - \boldsymbol{\eta}^{\top} \mathbf{h}$$

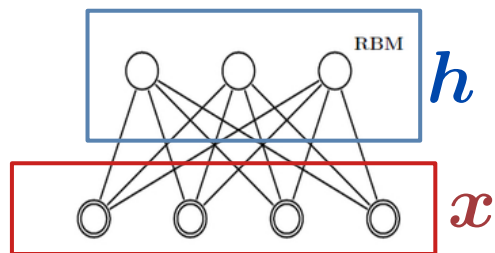
$$\boldsymbol{\theta} = \{W, \boldsymbol{\zeta}, \boldsymbol{\eta}\}$$

$$\begin{aligned} p_{\boldsymbol{\theta}}(\mathbf{x}) &= \frac{e^{-E_{\boldsymbol{\theta}}(\mathbf{x})}}{Z_{\boldsymbol{\theta}}} = \frac{\sum_{\mathbf{h}} e^{-\mathcal{E}_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{h})}}{Z_{\boldsymbol{\theta}}} = \frac{e^{\sum_i x_i \zeta_i}}{Z_{\boldsymbol{\theta}}} \prod_{a=1}^{N_h} \sum_{h_a=0}^1 e^{\sum_i x_i W_{ia} h_a + \eta_a h_a} \\ &= \frac{e^{\sum_i x_i \zeta_i}}{Z_{\boldsymbol{\theta}}} \prod_{a=1}^{N_h} \left(1 + e^{\sum_i x_i W_{ia} + \eta_a} \right) \end{aligned}$$

Models with hidden variables

- Restricted Boltzmann Machine

- Smolensky, P. (1986). *Information processing in dynamical systems: Foundations of harmony theory.*



visible variables
 $\mathcal{E}(\mathbf{x}, \mathbf{h}; \boldsymbol{\theta})$
Latent variables
Encode correlations

$$\mathcal{E}_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{h}) = -\mathbf{x}^{\top} \mathbf{W} \mathbf{h} - \boldsymbol{\zeta}^{\top} \mathbf{x} - \boldsymbol{\eta}^{\top} \mathbf{h}$$

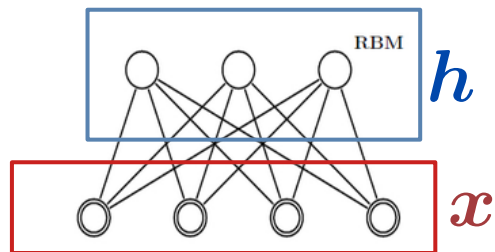
$$\boldsymbol{\theta} = \{\mathbf{W}, \boldsymbol{\zeta}, \boldsymbol{\eta}\}$$

$$\begin{aligned} p_{\boldsymbol{\theta}}(\mathbf{x}) &= \frac{e^{-E_{\boldsymbol{\theta}}(\mathbf{x})}}{Z_{\boldsymbol{\theta}}} = \frac{\sum_{\mathbf{h}} e^{-\mathcal{E}_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{h})}}{Z_{\boldsymbol{\theta}}} = \frac{e^{\sum_i x_i \zeta_i}}{Z_{\boldsymbol{\theta}}} \prod_{a=1}^{N_h} \sum_{h_a=0}^1 e^{\sum_i x_i W_{ia} h_a + \eta_a h_a} \\ &= \frac{e^{\sum_i x_i \zeta_i}}{Z_{\boldsymbol{\theta}}} \prod_{a=1}^{N_h} \left(1 + e^{\sum_i x_i W_{ia} + \eta_a} \right) \Rightarrow E_{\boldsymbol{\theta}}(\mathbf{x}) = -\sum_i x_i \zeta_i - \sum_{a=1}^{N_h} \log \left(1 + e^{\sum_i x_i W_{ia} + \eta_a} \right) \end{aligned}$$

Models with hidden variables

- Restricted Boltzmann Machine

- Smolensky, P. (1986). *Information processing in dynamical systems: Foundations of harmony theory.*



$\mathcal{E}(\mathbf{x}, \mathbf{h}; \theta)$

visible variables

Latent variables
Encode correlations

$$\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h}) = -\mathbf{x}^{\top} W \mathbf{h} - \zeta^{\top} \mathbf{x} - \eta^{\top} \mathbf{h}$$

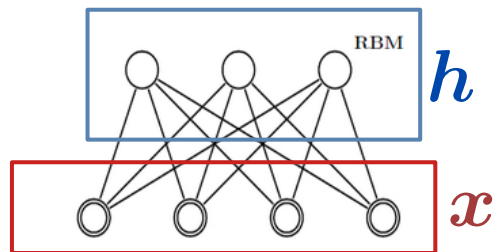
$$\theta = \{W, \zeta, \eta\}$$

$$\begin{aligned} p_{\theta}(\mathbf{x}) &= \frac{e^{-E_{\theta}(\mathbf{x})}}{Z_{\theta}} = \frac{\sum_{\mathbf{h}} e^{-\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h})}}{Z_{\theta}} = \frac{e^{\sum_i x_i \zeta_i}}{Z_{\theta}} \prod_{a=1}^{N_h} \sum_{h_a=0}^1 e^{\sum_i x_i W_{ia} h_a + \eta_a h_a} \\ &= \frac{e^{\sum_i x_i \zeta_i}}{Z_{\theta}} \prod_{a=1}^{N_h} \left(1 + e^{\sum_i x_i W_{ia} + \eta_a} \right) \Rightarrow E_{\theta}(\mathbf{x}) = -\sum_i x_i \zeta_i - \sum_{a=1}^{N_h} \log \left(1 + e^{\sum_i x_i W_{ia} + \eta_a} \right) \\ &\Rightarrow E_{\theta}(\mathbf{x}) = -\sum_i h_i x_i - \sum_{ij} J_{ij}^{(2)} x_i x_j - \sum_{ijk} J_{ijk}^{(3)} x_i x_j x_k - \sum_{ijkl} J_{ijkl}^{(4)} x_i x_j x_k x_l + \dots \end{aligned}$$

Models with hidden variables

- Restricted Boltzmann Machine

- Smolensky, P. (1986). *Information processing in dynamical systems: Foundations of harmony theory.*



visible variables
 $\mathcal{E}(\mathbf{x}, \mathbf{h}; \boldsymbol{\theta})$
 Latent variables
 Encode correlations

$$\mathcal{E}_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{h}) = -\mathbf{x}^{\top} W \mathbf{h} - \boldsymbol{\zeta}^{\top} \mathbf{x} - \boldsymbol{\eta}^{\top} \mathbf{h} \quad \boldsymbol{\theta} = \{W, \boldsymbol{\zeta}, \boldsymbol{\eta}\}$$

The marginal energy for the RBM encode high order interactions! → **Universal approximator**

Le Roux and Bengio. Neural computation (2008)

$$p_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{e^{-E_{\boldsymbol{\theta}}(\mathbf{x})}}{Z_{\boldsymbol{\theta}}} = \frac{\sum_{\mathbf{h}} e^{-\mathbf{x}^{\top} W \mathbf{h} - \boldsymbol{\zeta}^{\top} \mathbf{x} - \boldsymbol{\eta}^{\top} \mathbf{h}}}{Z_{\boldsymbol{\theta}}}$$

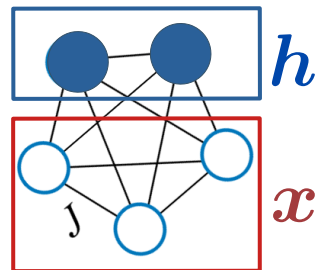
$$= \frac{e^{\sum_i x_i \zeta_i}}{Z_{\boldsymbol{\theta}}} \prod_{a=1}^{N_h} \left(1 + e^{\sum_i x_i w_{ia} + \eta_a} \right) \Rightarrow E_{\boldsymbol{\theta}}(\mathbf{x}) = -\sum_i x_i \zeta_i - \sum_{a=1} \log \left(1 + e^{\sum_i x_i w_{ia} + \eta_a} \right)$$

$$\Rightarrow E_{\boldsymbol{\theta}}(\mathbf{x}) = -\sum_i h_i x_i - \sum_{ij} J_{ij}^{(2)} x_i x_j - \sum_{ijk} J_{ijk}^{(3)} x_i x_j x_k - \sum_{ijkl} J_{ijkl}^{(4)} x_i x_j x_k x_l + \dots$$

Models with hidden variables

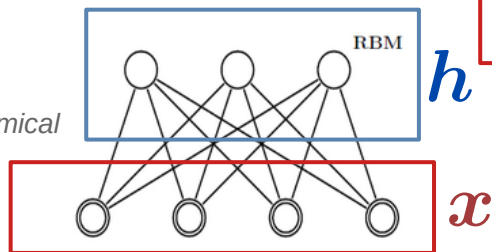
- Boltzmann Machines (Ising/Hopfield/Potts models)

- Ackley, D. H., Hinton, G. E., & Sejnowski, T. J. (1985). *A learning algorithm for Boltzmann machines*. *Cognitive science*, 9(1), 147-169.



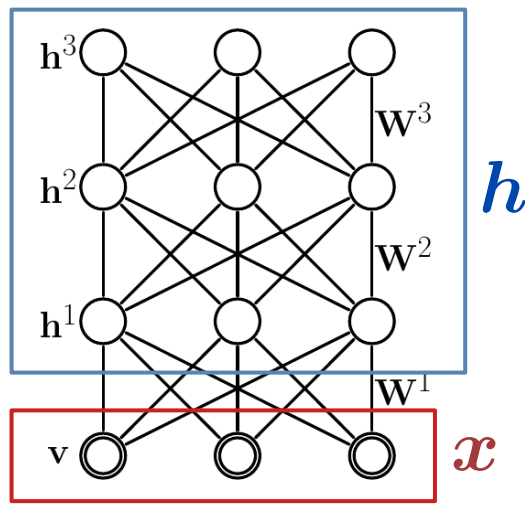
- Restricted Boltzmann Machine

- Smolensky, P. (1986). *Information processing in dynamical systems: Foundations of harmony theory*.



- Deep Boltzmann Machines

- Ruslan Salakhutdinov, Geoffrey Hinton (2009) *Deep Boltzmann Machines*.
- Bengio, Y. (2009). *Learning deep architectures for AI*.



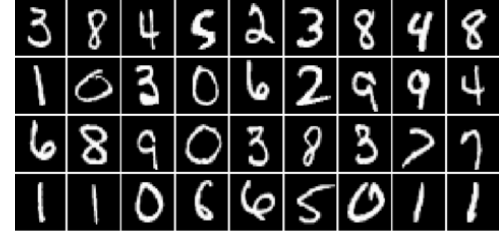
visible variables
 $\mathcal{E}(\mathbf{x}, \mathbf{h}; \theta)$
Latent variables
Encode correlations

$$p_{\theta}(\mathbf{x}) = \frac{\sum_{\mathbf{h}} e^{-\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h})}}{Z_{\theta}}$$
$$= \frac{e^{-E_{\theta}(\mathbf{x})}}{Z_{\theta}} \sim p_{\mathcal{D}}(\mathbf{x})$$



Training procedure

Training procedure



Goal of the training:

$$\text{Empirical } p_{\mathcal{D}}(\mathbf{x}) \sim \text{Model } p_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{e^{-E_{\boldsymbol{\theta}}(\mathbf{x})}}{Z_{\boldsymbol{\theta}}}$$

$$p_{\mathcal{D}}(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M \delta(\mathbf{x} - \mathbf{x}^{(m)})$$

Training procedure

Goal of the training:

$$p_{\mathcal{D}}(\mathbf{x}) \sim p_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{e^{-E_{\boldsymbol{\theta}}(\mathbf{x})}}{Z_{\boldsymbol{\theta}}}$$

Minimize
Kullback-Leibler (KL)
divergence

$$\begin{aligned} D_{\text{KL}}(p_{\mathcal{D}} || p_{\boldsymbol{\theta}}) &= \int d\mathbf{x} p_{\mathcal{D}}(\mathbf{x}) \log \frac{p_{\mathcal{D}}(\mathbf{x})}{p_{\boldsymbol{\theta}}(\mathbf{x})} \\ &= \underbrace{\int d\mathbf{x} p_{\mathcal{D}}(\mathbf{x}) \log p_{\mathcal{D}}(\mathbf{x})}_{\text{Constant}} - \int d\mathbf{x} p_{\mathcal{D}}(\mathbf{x}) \log p_{\boldsymbol{\theta}}(\mathbf{x}) \end{aligned}$$

Training procedure

$$p_{\mathcal{D}}(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M \delta(\mathbf{x} - \mathbf{x}^{(m)})$$

Goal of the training:

$$p_{\mathcal{D}}(\mathbf{x}) \sim p_{\theta}(\mathbf{x}) = \frac{e^{-E_{\theta}(\mathbf{x})}}{Z_{\theta}}$$

Minimize
Kullback-Leibler (KL)
divergence

$$\begin{aligned} D_{\text{KL}}(p_{\mathcal{D}} || p_{\theta}) &= \int d\mathbf{x} p_{\mathcal{D}}(\mathbf{x}) \log \frac{p_{\mathcal{D}}(\mathbf{x})}{p_{\theta}(\mathbf{x})} \\ &= \underbrace{\int d\mathbf{x} p_{\mathcal{D}}(\mathbf{x}) \log p_{\mathcal{D}}(\mathbf{x})}_{\text{Constant}} - \underbrace{\int d\mathbf{x} p_{\mathcal{D}}(\mathbf{x}) \log p_{\theta}(\mathbf{x})}_{\text{log-likelihood}} \end{aligned}$$

$$-\frac{1}{M} \sum_{m=1}^M \log p_{\theta}(\mathbf{x}^{(m)}) = -\frac{1}{M} \log \prod_{m=1}^M p_{\theta}(\mathbf{x}^{(m)}) = -\frac{1}{M} \log L(\mathcal{D} | \theta)$$

Training procedure

$$p_{\mathcal{D}}(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M \delta(\mathbf{x} - \mathbf{x}^{(m)})$$

Goal of the training:

$$p_{\mathcal{D}}(\mathbf{x}) \sim p_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{e^{-E_{\boldsymbol{\theta}}(\mathbf{x})}}{Z_{\boldsymbol{\theta}}}$$

Minimize

Kullback-Leibler
divergence

$$D_{\text{KL}}(p_{\mathcal{D}} || p_{\boldsymbol{\theta}}) \iff$$

Maximize

The log-
likelihood

$$\log L(\mathcal{D} | \boldsymbol{\theta}) \equiv \mathcal{L}(\mathcal{D} | \boldsymbol{\theta})$$

$$= \int d\mathbf{x} p_{\mathcal{D}}(\mathbf{x}) \log p_{\mathcal{D}}(\mathbf{x}) - \int d\mathbf{x} p_{\mathcal{D}}(\mathbf{x}) \log p_{\boldsymbol{\theta}}(\mathbf{x})$$

Constant

$$-\frac{1}{M} \sum_{m=1}^M \log p_{\boldsymbol{\theta}}(\mathbf{x}^{(m)}) = -\frac{1}{M} \log \prod_{m=1}^M p_{\boldsymbol{\theta}}(\mathbf{x}^{(m)}) = -\frac{1}{M} \log L(\mathcal{D} | \boldsymbol{\theta})$$

Training procedure

$$p_{\mathcal{D}}(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M \delta(\mathbf{x} - \mathbf{x}_d^{(m)})$$

Goal of the training:

$$p_{\mathcal{D}}(\mathbf{x}) \sim p_{\theta}(\mathbf{x}) = \frac{e^{-E_{\theta}(\mathbf{x})}}{Z_{\theta}}$$

Minimize

Kullback-Leibler divergence

$$D_{\text{KL}}(p_{\mathcal{D}} || p_{\theta})$$



Maximize

The log-likelihood

$$\log L(\mathcal{D} | \theta) \equiv \mathcal{L}(\mathcal{D} | \theta)$$

Kullback-Leibler divergence

$$= \int d\mathbf{x} p_{\mathcal{D}}(\mathbf{x}) \log p_{\mathcal{D}}(\mathbf{x}) - \int d\mathbf{x} p_{\mathcal{D}}(\mathbf{x}) \log p_{\theta}(\mathbf{x})$$

Recall Bayes-Theorem

$$\underbrace{p(\mathcal{D} | \theta)}_{\text{likelihood}} p(\theta) = p(\theta | \mathcal{D}) \underbrace{p(\mathcal{D})}_{\text{posterior}}$$

constant

$$-\frac{1}{M} \log \prod_{m=1}^M p_{\theta}(\mathbf{x}^{(m)}) = -\frac{1}{M} \log L(\mathcal{D} | \theta)$$

Log-likelihood maximization

$$\mathcal{L}(\mathcal{D}|\boldsymbol{\theta}) = \sum_{m=1}^M \log p_{\boldsymbol{\theta}}(\mathbf{x} = \mathbf{x}^{(m)})$$

$$p_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{e^{-E_{\boldsymbol{\theta}}(\mathbf{x})}}{Z_{\boldsymbol{\theta}}}$$

Log-likelihood maximization

$$\mathcal{L}(\mathcal{D}|\boldsymbol{\theta}) = \sum_{m=1}^M \log p_{\boldsymbol{\theta}}(\mathbf{x} = \mathbf{x}^{(m)})$$

$$p_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{e^{-E_{\boldsymbol{\theta}}(\mathbf{x})}}{Z_{\boldsymbol{\theta}}}$$

Partition function

$$\mathcal{L}(\mathcal{D}|\boldsymbol{\theta}) = \langle \log p_{\boldsymbol{\theta}}(\mathbf{x}) \rangle_{p_{\mathcal{D}}} = \langle -E_{\boldsymbol{\theta}}(\mathbf{x}) \rangle_{p_{\mathcal{D}}} - \log Z_{\boldsymbol{\theta}}$$

$$Z_{\boldsymbol{\theta}} = \sum_{\{\mathbf{x}\}} e^{-E_{\boldsymbol{\theta}}(\mathbf{x})}$$

If x_i binary $\rightarrow 2^N$

Intractable

Log-likelihood maximization

$$\mathcal{L}(\mathcal{D}|\boldsymbol{\theta}) = \sum_{m=1}^M \log p_{\boldsymbol{\theta}}(\mathbf{x} = \mathbf{x}^{(m)})$$

$$p_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{e^{-E_{\boldsymbol{\theta}}(\mathbf{x})}}{Z_{\boldsymbol{\theta}}}$$

$$\mathcal{L}(\mathcal{D}|\boldsymbol{\theta}) = \langle \log p_{\boldsymbol{\theta}}(\mathbf{x}) \rangle_{p_{\mathcal{D}}} = \langle -E_{\boldsymbol{\theta}}(\mathbf{x}) \rangle_{p_{\mathcal{D}}} - \log Z_{\boldsymbol{\theta}}$$

Partition function

$$Z_{\boldsymbol{\theta}} = \sum_{\{\mathbf{x}\}} e^{-E_{\boldsymbol{\theta}}(\mathbf{x})}$$

If x_i binary $\rightarrow 2^N$

Intractable

(Stochastic) gradient ascent

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}$$

$$\theta_i^{(t+1)} \leftarrow \theta_i^t + \gamma \left. \frac{\partial \mathcal{L}}{\partial \theta_i} \right|_{\theta = \theta_i^{(t)}}$$

Log-likelihood maximization

$$\frac{\partial \mathcal{L}}{\partial \theta_i} = \left\langle -\frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\mathcal{D}}} - \frac{\partial \log Z}{\partial \theta_i}$$

Partition function

$$\mathcal{L}(\mathcal{D}|\boldsymbol{\theta}) = \langle \log p_{\boldsymbol{\theta}}(\mathbf{x}) \rangle_{p_{\mathcal{D}}} = \langle -E_{\boldsymbol{\theta}}(\mathbf{x}) \rangle_{p_{\mathcal{D}}} - \log Z_{\boldsymbol{\theta}}$$

(Stochastic) gradient ascent

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}$$

$$\theta_i^{(t+1)} \leftarrow \theta_i^t + \gamma \frac{\partial \mathcal{L}}{\partial \theta_i} \Big|_{\theta = \theta_i^{(t)}}$$

Log-likelihood maximization

$$\frac{\partial \mathcal{L}}{\partial \theta_i} = \left\langle -\frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\mathcal{D}}} - \frac{\partial \log Z}{\partial \theta_i}$$

$$\begin{aligned} \frac{\partial \log Z}{\partial \theta_i} &= \sum_{\{\mathbf{x}\}} \frac{e^{-E(\mathbf{x})}}{Z} \frac{\partial E(\mathbf{x})}{\partial \theta_i} \\ &= \left\langle \frac{\partial E(\mathbf{x})}{\partial \theta_i} \right\rangle_{p_{\theta}(\mathbf{x})} \end{aligned}$$

$$\mathcal{L}(\mathcal{D}|\theta) = \langle \log p_{\theta}(\mathbf{x}) \rangle_{p_{\mathcal{D}}} = \langle -E_{\theta}(\mathbf{x}) \rangle_{p_{\mathcal{D}}} - \log Z_{\theta}$$

(Stochastic) gradient ascent

$$\nabla_{\theta} \mathcal{L}$$

$$\theta_i^{(t+1)} \leftarrow \theta_i^t + \gamma \left. \frac{\partial \mathcal{L}}{\partial \theta_i} \right|_{\theta=\theta_i^{(t)}}$$

$$\nabla \mathcal{L}_{\theta} = \underbrace{\langle -\nabla E_{\theta} \rangle_{p_{\mathcal{D}}}}_{\text{data}} - \underbrace{\langle -\nabla E_{\theta} \rangle_{p_{\theta}}}_{\text{model}}$$

Log-likelihood maximization

$$\frac{\partial \mathcal{L}}{\partial \theta_i} = \left\langle -\frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\mathcal{D}}} - \frac{\partial \log Z}{\partial \theta_i}$$

$$\begin{aligned} \frac{\partial \log Z}{\partial \theta_i} &= \sum_{\{\mathbf{x}\}} \frac{e^{-E(\mathbf{x})}}{Z} \frac{\partial E(\mathbf{x})}{\partial \theta_i} \\ &= \left\langle \frac{\partial E(\mathbf{x})}{\partial \theta_i} \right\rangle_{p_{\theta}(\mathbf{x})} \end{aligned}$$

$$\nabla E_{\theta}^{\mathcal{D}}$$

$$p_{\mathcal{D}}(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M \delta(\mathbf{x} - \mathbf{x}_d^{(m)})$$

(Stochastic) gradient ascent

$$\nabla_{\theta} \mathcal{L}$$

$$\theta_i^{(t+1)} \leftarrow \theta_i^t + \gamma \left. \frac{\partial \mathcal{L}}{\partial \theta_i} \right|_{\theta = \theta_i^{(t)}}$$

$$\nabla \mathcal{L}_{\theta} = \underbrace{\left\langle -\nabla E_{\theta} \right\rangle_{p_{\mathcal{D}}}}_{\text{data}} - \underbrace{\left\langle -\nabla E_{\theta} \right\rangle_{p_{\theta}}}_{\text{model}}$$

Easy !

Log-likelihood maximization

$$\frac{\partial \mathcal{L}}{\partial \theta_i} = \left\langle -\frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\mathcal{D}}} - \frac{\partial \log Z}{\partial \theta_i}$$

$$\begin{aligned} \frac{\partial \log Z}{\partial \theta_i} &= \sum_{\{\mathbf{x}\}} \frac{e^{-E(\mathbf{x})}}{Z} \frac{\partial E(\mathbf{x})}{\partial \theta_i} \\ &= \left\langle \frac{\partial E(\mathbf{x})}{\partial \theta_i} \right\rangle_{p_{\theta}} \end{aligned}$$

$$p_{\theta}(\mathbf{x}) = \frac{e^{-E_{\theta}(\mathbf{x})}}{Z_{\theta}}$$

$$p_{\mathcal{D}}(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M \delta(\mathbf{x} - \mathbf{x}_d^{(m)})$$

MCMC
sampling

(Stochastic) gradient ascent

$$\nabla_{\theta} \mathcal{L}$$

$$\theta_i^{(t+1)} \leftarrow \theta_i^t + \gamma \frac{\partial \mathcal{L}}{\partial \theta_i} \Big|_{\theta = \theta_i^{(t)}}$$

$$\nabla \mathcal{L}_{\theta} = \underbrace{\left\langle -\nabla E_{\theta} \right\rangle_{p_{\mathcal{D}}}}_{\text{data}} - \underbrace{\left\langle -\nabla E_{\theta} \right\rangle_{p_{\theta}}}_{\text{model}}$$

Easy !

Hard !

Log-likelihood

Every time we want to update the parameters

$$\mathbf{x}_{\text{gen}}^{(m)} \quad m = 1, \dots, n_{\text{chains}}$$

$$\mathbf{X}_{\text{gen}} \sim P_{\theta} \quad \text{Via a Markov Chain Monte Carlo process}$$

$$\langle -\nabla E_{\theta} \rangle_{p_{\theta}} \approx \frac{1}{n_{\text{chains}}} \sum_{m=1}^{n_{\text{chains}}} \nabla E(\mathbf{x}_{\text{gen}}^{(m)})$$

$$p_{\theta}(\mathbf{x}) = \frac{e^{-E_{\theta}(\mathbf{x})}}{Z_{\theta}}$$

$$p_{\mathcal{D}}(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M \delta(\mathbf{x} - \mathbf{x}_d^{(m)})$$

MCMC sampling

(Stochastic) gradient ascent

$$\nabla_{\theta} \mathcal{L}$$

$$\theta_i^{(t+1)} \leftarrow \theta_i^t + \gamma \left. \frac{\partial \mathcal{L}}{\partial \theta_i} \right|_{\theta = \theta_i^{(t)}}$$

$$\nabla \mathcal{L}_{\theta} = \underbrace{\langle -\nabla E_{\theta} \rangle_{p_{\mathcal{D}}}}_{\text{data}} - \underbrace{\langle -\nabla E_{\theta} \rangle_{p_{\theta}}}_{\text{model}}$$

Easy !

Hard !

Log-likelihood

Every time we want to update the parameters

$$\mathbf{x}_{\text{gen}}^{(m)} \quad m = 1, \dots, n_{\text{chains}}$$

$$\mathbf{X}_{\text{gen}} \sim P_{\theta} \quad \text{Via a Markov Chain Monte Carlo process}$$

$$\langle -\nabla E_{\theta} \rangle_{p_{\theta}} \approx \frac{1}{n_{\text{chains}}} \sum_{m=1}^{n_{\text{chains}}} \nabla E(\mathbf{x}_{\text{gen}}^{(m)})$$

Origin of all the difficulties ! → 3rd lecture

$$p_{\theta}(\mathbf{x}) = \frac{e^{-E_{\theta}(\mathbf{x})}}{Z_{\theta}}$$

$$p_{\mathcal{D}}(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M \delta(\mathbf{x} - \mathbf{x}_d^{(m)})$$

MCMC sampling

(Stochastic) gradient ascent

$$\nabla_{\theta} \mathcal{L}$$

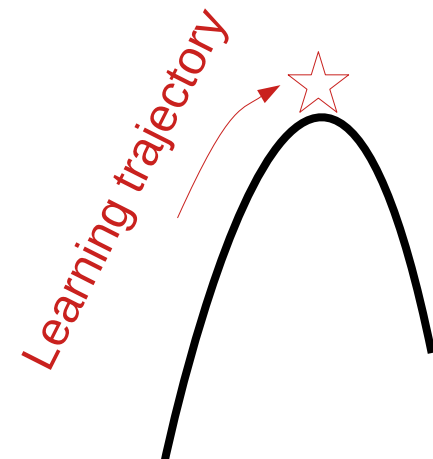
$$\theta_i^{(t+1)} \leftarrow \theta_i^t + \gamma \frac{\partial \mathcal{L}}{\partial \theta_i} \Big|_{\theta = \theta_i^{(t)}}$$

$$\nabla \mathcal{L}_{\theta} = \underbrace{\langle -\nabla E_{\theta} \rangle_{p_{\mathcal{D}}}}_{\text{data}} - \underbrace{\langle -\nabla E_{\theta} \rangle_{p_{\theta}}}_{\text{model}}$$

Easy !

Hard !

On the gradient ascent



$$\theta(t + t) \leftarrow \theta(t) + \gamma \nabla \mathcal{L}(t)$$

Update rule:

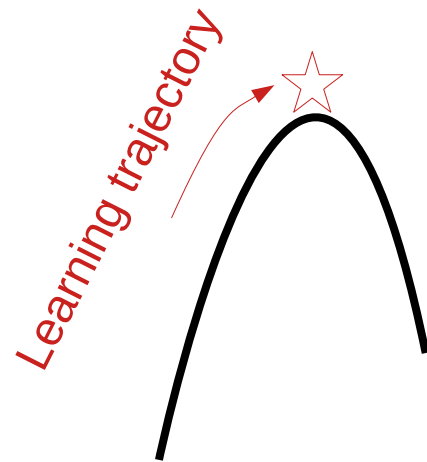
$$\nabla \mathcal{L}_{\theta} = \langle -\nabla E_{\theta} \rangle_{p_{\mathcal{D}}} - \langle -\nabla E_{\theta} \rangle_{p_{\theta}}$$

On the gradient ascent

★ Fixed point : $\nabla \mathcal{L}_{\theta} = \mathbf{0}$

$$\left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\mathcal{D}}} = \left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\theta}} \quad \forall \theta_i$$

Moment matching statistics



$$\theta(t + t) \leftarrow \theta(t) + \gamma \nabla \mathcal{L}(t)$$

Update rule:

$$\nabla \mathcal{L}_{\theta} = \langle -\nabla E_{\theta} \rangle_{p_{\mathcal{D}}} - \langle -\nabla E_{\theta} \rangle_{p_{\theta}}$$

On the gradient ascent

$$f_{\theta_i}(\mathbf{x}, \boldsymbol{\theta}) \equiv \frac{\partial E_{\boldsymbol{\theta}}(\mathbf{x})}{\partial \theta_i}$$

★ Fixed point : $\nabla \mathcal{L}_{\boldsymbol{\theta}} = \mathbf{0}$

$$\left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\mathcal{D}}} = \left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\boldsymbol{\theta}}} \quad \forall \theta_i$$

Moment matching statistics

Hessian matrix

$$H_{ij}(\boldsymbol{\theta}) \equiv \frac{\partial^2 \mathcal{L}(\boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} = \left\langle \frac{\partial f_{\theta_j}(\mathbf{x})}{\partial \theta_i} \right\rangle_{p_{\boldsymbol{\theta}}} - \left\langle \frac{\partial f_{\theta_j}(\mathbf{x})}{\partial \theta_i} \right\rangle_{p_{\mathcal{D}}} \\ - \langle f_{\theta_i}(\mathbf{x}, \boldsymbol{\theta}) f_{\theta_j}(\mathbf{x}, \boldsymbol{\theta}) \rangle_{p_{\boldsymbol{\theta}}} + \langle f_{\theta_i}(\mathbf{x}, \boldsymbol{\theta}) \rangle_{p_{\boldsymbol{\theta}}} \langle f_{\theta_j}(\mathbf{x}, \boldsymbol{\theta}) \rangle_{p_{\boldsymbol{\theta}}}$$

Update rule:

$$\nabla \mathcal{L}_{\boldsymbol{\theta}} = \langle -\nabla E_{\boldsymbol{\theta}} \rangle_{p_{\mathcal{D}}} - \langle -\nabla E_{\boldsymbol{\theta}} \rangle_{p_{\boldsymbol{\theta}}}$$

On the gradient ascent

$$f_{\theta_i}(\mathbf{x}, \boldsymbol{\theta}) \equiv \frac{\partial E_{\boldsymbol{\theta}}(\mathbf{x})}{\partial \theta_i}$$

★ Fixed point : $\nabla \mathcal{L}_{\boldsymbol{\theta}} = \mathbf{0}$

$$\left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\mathcal{D}}} = \left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\boldsymbol{\theta}}} \quad \forall \theta_i$$

Moment matching statistics

Hessian matrix

$$H_{ij}(\boldsymbol{\theta}) \equiv \frac{\partial^2 \mathcal{L}(\boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} = \left\langle \frac{\partial f_{\theta_j}(\mathbf{x})}{\partial \theta_i} \right\rangle_{p_{\boldsymbol{\theta}}} \overset{\text{if}}{\cancel{\left\langle \frac{\partial f_{\theta_j}(\mathbf{x})}{\partial \theta_i} \right\rangle_{p_{\mathcal{D}}}}} - \left\langle f_{\theta_i}(\mathbf{x}, \boldsymbol{\theta}) f_{\theta_j}(\mathbf{x}, \boldsymbol{\theta}) \right\rangle_{p_{\boldsymbol{\theta}}} + \left\langle f_{\theta_i}(\mathbf{x}, \boldsymbol{\theta}) \right\rangle_{p_{\boldsymbol{\theta}}} \left\langle f_{\theta_j}(\mathbf{x}, \boldsymbol{\theta}) \right\rangle_{p_{\boldsymbol{\theta}}}$$

Semi negative definite \Rightarrow convex

Update rule:

$$\nabla \mathcal{L}_{\boldsymbol{\theta}} = \left\langle -\nabla E_{\boldsymbol{\theta}} \right\rangle_{p_{\mathcal{D}}} - \left\langle -\nabla E_{\boldsymbol{\theta}} \right\rangle_{p_{\boldsymbol{\theta}}}$$

Example 1: Boltzmann Machine

★ Fixed point : $\nabla \mathcal{L}_{\theta} = \mathbf{0}$

$$\left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\mathcal{D}}} = \left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\theta}} \quad \forall \theta_i$$

Moment matching statistics

$$\frac{\partial \mathcal{L}}{\partial J_{ij}} = \langle S_i S_j \rangle_{p_{\mathcal{D}}} - \langle S_i S_j \rangle_{p_{\theta}}$$

$$\frac{\partial \mathcal{L}}{\partial h_i} = \langle S_i \rangle_{p_{\mathcal{D}}} - \langle S_i \rangle_{p_{\theta}}$$

Ising-like model

$$E_{J, \mathbf{h}}(\mathbf{S}) = - \sum_{ij} J_{ij} S_i S_j - \sum_i h_i S_i$$

$$\frac{\partial E}{\partial J_{ij}} = -S_i S_j \quad \frac{\partial E}{\partial h_i} = -S_i$$

Example 1: Boltzmann Machine

★ Fixed point : $\nabla \mathcal{L}_{\theta} = \mathbf{0}$

$$\left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\mathcal{D}}} = \left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\theta}} \quad \forall \theta_i$$

Moment matching statistics

Ising-like model

$$E_{J, \mathbf{h}}(\mathbf{S}) = - \sum_{ij} J_{ij} S_i S_j - \sum_i h_i S_i$$

Example 1: Boltzmann Machine

★ Fixed point : $\nabla \mathcal{L}_{\theta} = \mathbf{0}$

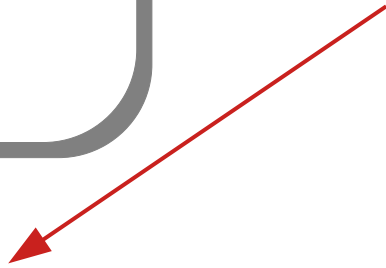
$$\left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\mathcal{D}}} = \left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\theta}} \quad \forall \theta_i$$

Moment matching statistics

Ising-like model

$$E_{J, \mathbf{h}}(\mathbf{S}) = - \sum_{ij} J_{ij} S_i S_j - \sum_i h_i S_i$$

$$\frac{\partial E}{\partial J_{ij}} = -S_i S_j \quad \frac{\partial E}{\partial h_i} = -S_i$$



$$\frac{\partial \mathcal{L}}{\partial J_{ij}} = \langle S_i S_j \rangle_{p_{\mathcal{D}}} - \langle S_i S_j \rangle_{p_{\theta}}$$

$$\frac{\partial \mathcal{L}}{\partial h_i} = \langle S_i \rangle_{p_{\mathcal{D}}} - \langle S_i \rangle_{p_{\theta}}$$

Example 1: Boltzmann Machine

★ Fixed point : $\nabla \mathcal{L}_{\theta} = \mathbf{0}$

$$\left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\mathcal{D}}} = \left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\theta}} \quad \forall \theta_i$$

Moment matching statistics

$$\frac{\partial \mathcal{L}}{\partial J_{ij}} = \langle S_i S_j \rangle_{p_{\mathcal{D}}} - \langle S_i S_j \rangle_{p_{\theta}}$$

$$\frac{\partial \mathcal{L}}{\partial h_i} = \langle S_i \rangle_{p_{\mathcal{D}}} - \langle S_i \rangle_{p_{\theta}}$$

Ising-like model

$$E_{J,h}(\mathbf{S}) = - \sum_{ij} J_{ij} S_i S_j - \sum_i h_i S_i$$

$$\frac{\partial E}{\partial J_{ij}} = -S_i S_j \quad \frac{\partial E}{\partial h_i} = -S_i$$

We can encode the covariance matrix of the data but nothing beyond that!

Solution is unique !

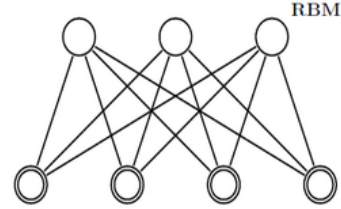
Fixed point

$$\langle S_i S_j \rangle_{p_{J,h}} = \langle S_i S_j \rangle_{p_{\mathcal{D}}}$$

$$\langle S_i \rangle_{p_{J,h}} = \langle S_i \rangle_{p_{\mathcal{D}}}$$

Example 2: Restricted Boltzmann Machine

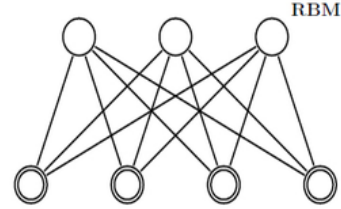
$$\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h}) = -\mathbf{x}^{\top} \mathbf{W} \mathbf{h} - \boldsymbol{\zeta}^{\top} \mathbf{x} - \boldsymbol{\eta}^{\top} \mathbf{h} \quad p_{\theta}(\mathbf{x}) = \frac{\sum_{\mathbf{h}} e^{-\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h})}}{Z_{\theta}}$$



$$h_a = \{0, 1\} \Rightarrow E_{\theta}(\mathbf{x}) = -\sum_i x_i \zeta_i - \sum_{a=1}^{N_h} \log \left(1 + e^{\sum_i x_i W_{ia} + \eta_a} \right)$$

Example 2: Restricted Boltzmann Machine

$$\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h}) = -\mathbf{x}^{\top} \mathbf{W} \mathbf{h} - \boldsymbol{\zeta}^{\top} \mathbf{x} - \boldsymbol{\eta}^{\top} \mathbf{h} \quad p_{\theta}(\mathbf{x}) = \frac{\sum_{\mathbf{h}} e^{-\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h})}}{Z_{\theta}}$$



$$h_a = \{0, 1\} \Rightarrow E_{\theta}(\mathbf{x}) = -\sum_i x_i \zeta_i - \sum_{a=1}^{N_h} \log \left(1 + e^{\sum_i x_i W_{ia} + \eta_a} \right)$$

$$\frac{\partial E}{\partial W_{ia}} = -\sigma \left(\sum_a W_{ia} x_i + \eta_a \right) x_i$$

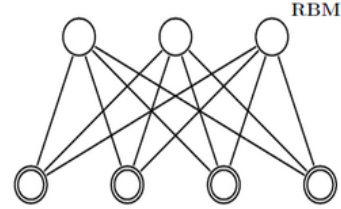
$$\frac{\partial E}{\partial \eta_a} = -\sigma \left(\sum_a W_{ia} x_i + \eta_a \right)$$

$$\frac{\partial E}{\partial \zeta_i} = -x_i$$

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \text{sigmoid}(x)$$

Example 2: Restricted Boltzmann Machine

$$\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h}) = -\mathbf{x}^{\top} \mathbf{W} \mathbf{h} - \boldsymbol{\zeta}^{\top} \mathbf{x} - \boldsymbol{\eta}^{\top} \mathbf{h} \quad p_{\theta}(\mathbf{x}) = \frac{\sum_{\mathbf{h}} e^{-\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h})}}{Z_{\theta}}$$



$$h_a = \{0, 1\} \Rightarrow E_{\theta}(\mathbf{x}) = -\sum_i x_i \zeta_i - \sum_{a=1}^{N_h} \log \left(1 + e^{\sum_i x_i W_{ia} + \eta_a} \right)$$

$$\frac{\partial E}{\partial W_{ia}} = -\sigma \left(\sum_a W_{ia} x_i + \eta_a \right) x_i$$

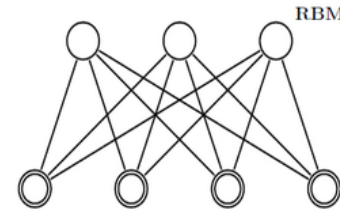
$$\frac{\partial E}{\partial \eta_a} = -\sigma \left(\sum_a W_{ia} x_i + \eta_a \right)$$

$$\frac{\partial E}{\partial \zeta_i} = -x_i$$

$$\begin{aligned} p(h_a = 1 | \mathbf{x}, \mathbf{h}_{-a}, \boldsymbol{\theta}) &= \frac{e^{\sum_i W_{ia} x_i + \eta_a}}{1 + e^{\sum_i W_{ia} x_i + \eta_a}} \\ &= \sigma \left(\sum_i W_{ia} x_i + \eta_a \right) = \langle h_a \rangle_{p_{\mathcal{E}}(\mathbf{h} | \mathbf{x})} \end{aligned}$$

Example 2: Restricted Boltzmann Machine

$$\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h}) = -\mathbf{x}^{\top} \mathbf{W} \mathbf{h} - \boldsymbol{\zeta}^{\top} \mathbf{x} - \boldsymbol{\eta}^{\top} \mathbf{h} \quad p_{\theta}(\mathbf{x}) = \frac{\sum_{\mathbf{h}} e^{-\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h})}}{Z_{\theta}}$$



$$h_a = \{0, 1\} \Rightarrow E_{\theta}(\mathbf{x}) = -\sum_i x_i \zeta_i - \sum_{a=1}^{N_h} \log \left(1 + e^{\sum_i x_i W_{ia} + \eta_a} \right)$$

$$\frac{\partial E}{\partial W_{ia}} = -\sigma \left(\sum_a W_{ia} x_i + \eta_a \right) x_i = -\langle h_a \rangle_{p_{\mathcal{E}}(\mathbf{h}|\mathbf{x})} x_i$$

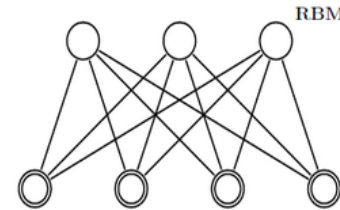
$$\frac{\partial E}{\partial \eta_a} = -\sigma \left(\sum_a W_{ia} x_i + \eta_a \right) = -\langle h_a \rangle_{p_{\mathcal{E}}(\mathbf{h}|\mathbf{x})}$$

$$\frac{\partial E}{\partial \zeta_i} = -x_i$$

$$p(h_a = 1 | \mathbf{x}, \mathbf{h}_{-a}, \boldsymbol{\theta}) = \frac{e^{\sum_i W_{ia} x_i + \eta_a}}{1 + e^{\sum_i W_{ia} x_i + \eta_a}} = \sigma \left(\sum_i W_{ia} x_i + \eta_a \right) = \langle h_a \rangle_{p_{\mathcal{E}}(\mathbf{h}|\mathbf{x})}$$

Example 2: Restricted Boltzmann Machine

$$\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h}) = -\mathbf{x}^{\top} W \mathbf{h} - \zeta^{\top} \mathbf{x} - \eta^{\top} \mathbf{h} \quad p_{\theta}(\mathbf{x}) = \frac{\sum_{\mathbf{h}} e^{-\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h})}}{Z_{\theta}}$$



$$h_a = \{0, 1\} \Rightarrow E_{\theta}(\mathbf{x}) = -\sum_i x_i \zeta_i - \sum_{a=1}^{N_h} \log \left(1 + e^{\sum_i x_i W_{ia} + \eta_a} \right)$$

$$p_{\mathcal{E}}(\mathbf{x}, \mathbf{h}) = p_{\mathcal{E}}(\mathbf{h}|\mathbf{x}) p_{\theta}(\mathbf{x})$$

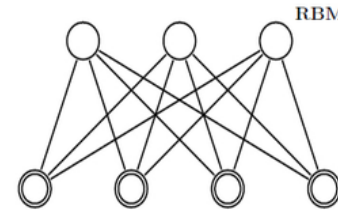
$$\frac{\partial \mathcal{L}}{\partial W_{ia}} = \left\langle x_i \langle h_a \rangle_{p_{\mathcal{E}}(\mathbf{h}|\mathbf{x})} \right\rangle_{p_{\mathcal{D}}} - \left\langle x_i \langle h_a \rangle_{p_{\mathcal{E}}(\mathbf{h}|\mathbf{x})} \right\rangle_{p_{\theta}}$$

$$\frac{\partial \mathcal{L}}{\partial \eta_a} = \left\langle \langle h_a \rangle_{p_{\mathcal{E}}(\mathbf{h}|\mathbf{x})} \right\rangle_{p_{\mathcal{D}}} - \left\langle \langle h_a \rangle_{p_{\mathcal{E}}(\mathbf{h}|\mathbf{x})} \right\rangle_{p_{\theta}}$$

$$\frac{\partial \mathcal{L}}{\partial \zeta_i} = \langle x_i \rangle_{p_{\mathcal{D}}} - \langle x_i \rangle_{p_{\theta}}$$

Example 2: Restricted Boltzmann Machine

$$\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h}) = -\mathbf{x}^{\top} W \mathbf{h} - \zeta^{\top} \mathbf{x} - \eta^{\top} \mathbf{h} \quad p_{\theta}(\mathbf{x}) = \frac{\sum_{\mathbf{h}} e^{-\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h})}}{Z_{\theta}}$$



$$h_a = \{0, 1\} \Rightarrow E_{\theta}(\mathbf{x}) = -\sum_i x_i \zeta_i - \sum_{a=1}^{N_h} \log \left(1 + e^{\sum_i x_i W_{ia} + \eta_a} \right)$$

$$p_{\mathcal{E}}(\mathbf{x}, \mathbf{h}) = p_{\mathcal{E}}(\mathbf{h}|\mathbf{x}) p_{\theta}(\mathbf{x})$$

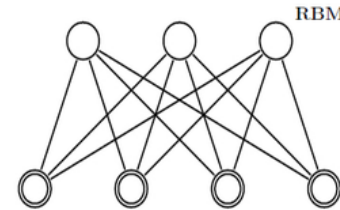
$$\frac{\partial \mathcal{L}}{\partial W_{ia}} = \left\langle x_i \langle h_a \rangle_{p_{\mathcal{E}}(\mathbf{h}|\mathbf{x})} \right\rangle_{p_{\mathcal{D}}} - \langle x_i h_a \rangle_{\mathcal{E}}$$

$$\frac{\partial \mathcal{L}}{\partial \eta_a} = \left\langle \langle h_a \rangle_{p_{\mathcal{E}}(\mathbf{h}|\mathbf{x})} \right\rangle_{p_{\mathcal{D}}} - \langle h_a \rangle_{\mathcal{E}}$$

$$\frac{\partial \mathcal{L}}{\partial \zeta_i} = \langle x_i \rangle_{p_{\mathcal{D}}} - \langle x_i \rangle_{\mathcal{E}}$$

Example 2: Restricted Boltzmann Machine

$$\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h}) = -\mathbf{x}^{\top} \mathbf{W} \mathbf{h} - \boldsymbol{\zeta}^{\top} \mathbf{x} - \boldsymbol{\eta}^{\top} \mathbf{h} \quad p_{\theta}(\mathbf{x}) = \frac{\sum_{\mathbf{h}} e^{-\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h})}}{Z_{\theta}}$$



$$h_a = \{0, 1\} \Rightarrow E_{\theta}(\mathbf{x}) = -\sum_i x_i \zeta_i - \sum_{a=1}^{N_h} \log \left(1 + e^{\sum_i x_i W_{ia} + \eta_a} \right)$$

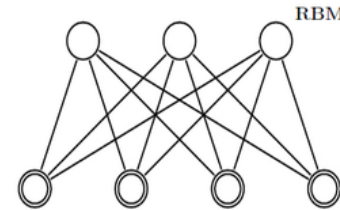
$$\frac{\partial \mathcal{L}}{\partial W_{ia}} = \langle x_i h_a \rangle_{p_{\mathcal{D}}} - \langle x_i h_a \rangle_{\mathcal{E}}$$

$$\frac{\partial \mathcal{L}}{\partial \eta_a} = \langle h_a \rangle_{p_{\mathcal{D}}} - \langle h_a \rangle_{\mathcal{E}}$$

$$\frac{\partial \mathcal{L}}{\partial \zeta_i} = \langle x_i \rangle_{p_{\mathcal{D}}} - \langle x_i \rangle_{\mathcal{E}}$$

Example 2: Restricted Boltzmann Machine

$$\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h}) = -\mathbf{x}^{\top} \mathbf{W} \mathbf{h} - \boldsymbol{\zeta}^{\top} \mathbf{x} - \boldsymbol{\eta}^{\top} \mathbf{h} \quad p_{\theta}(\mathbf{x}) = \frac{\sum_{\mathbf{h}} e^{-\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h})}}{Z_{\theta}}$$



$$h_a = \{0, 1\} \Rightarrow E_{\theta}(\mathbf{x}) = -\sum_i x_i \zeta_i - \sum_{a=1}^{N_h} \log \left(1 + e^{\sum_i x_i W_{ia} + \eta_a} \right)$$

Boltzmann machine:

$$\frac{\partial \mathcal{L}}{\partial J_{ij}} = \langle S_i S_j \rangle_{p_{\mathcal{D}}} - \langle S_i S_j \rangle_{p_{\theta}}$$

$$\frac{\partial \mathcal{L}}{\partial h_i} = \langle S_i \rangle_{p_{\mathcal{D}}} - \langle S_i \rangle_{p_{\theta}}$$

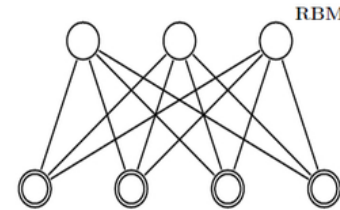
$$\frac{\partial \mathcal{L}}{\partial W_{ia}} = \langle x_i h_a \rangle_{p_{\mathcal{D}}} - \langle x_i h_a \rangle_{\mathcal{E}}$$

$$\frac{\partial \mathcal{L}}{\partial \eta_a} = \langle h_a \rangle_{p_{\mathcal{D}}} - \langle h_a \rangle_{\mathcal{E}}$$

$$\frac{\partial \mathcal{L}}{\partial \zeta_i} = \langle x_i \rangle_{p_{\mathcal{D}}} - \langle x_i \rangle_{\mathcal{E}}$$

Example 2: Restricted Boltzmann Machine

$$\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h}) = -\mathbf{x}^{\top} \mathbf{W} \mathbf{h} - \boldsymbol{\zeta}^{\top} \mathbf{x} - \boldsymbol{\eta}^{\top} \mathbf{h} \quad p_{\theta}(\mathbf{x}) = \frac{\sum_{\mathbf{h}} e^{-\mathcal{E}_{\theta}(\mathbf{x}, \mathbf{h})}}{Z_{\theta}}$$



$$h_a = \{0, 1\} \Rightarrow E_{\theta}(\mathbf{x}) = -\sum_i x_i \zeta_i - \sum_{a=1}^{N_h} \log \left(1 + e^{\sum_i x_i W_{ia} + \eta_a} \right)$$

Boltzmann machine:

$$\frac{\partial \mathcal{L}}{\partial J_{ij}} = \langle S_i S_j \rangle_{p_{\mathcal{D}}} - \langle S_i S_j \rangle_{p_{\theta}}$$

$$\frac{\partial \mathcal{L}}{\partial h_i} = \langle S_i \rangle_{p_{\mathcal{D}}} - \langle S_i \rangle_{p_{\theta}}$$

$$\frac{\partial \mathcal{L}}{\partial W_{ia}} = \langle x_i h_a \rangle_{p_{\mathcal{D}}} - \langle x_i h_a \rangle_{\mathcal{E}}$$

$$\frac{\partial \mathcal{L}}{\partial \eta_a} = \langle h_a \rangle_{p_{\mathcal{D}}} - \langle h_a \rangle_{\mathcal{E}}$$

$$\frac{\partial \mathcal{L}}{\partial \zeta_i} = \langle x_i \rangle_{p_{\mathcal{D}}} - \langle x_i \rangle_{\mathcal{E}}$$



Sample generation

Generating new samples

Empirical

Model

$$p_{\mathcal{D}}(\mathbf{x}) \sim \frac{e^{-E_{\theta}(\mathbf{x})}}{Z_{\theta}}$$

Dominated minimum
free-energy
configurations

$$\{\mathbf{x}\}_{\text{eq},\theta} \sim \mathcal{D}$$

Generating new samples

Empirical

Model

$$p_{\mathcal{D}}(\mathbf{x}) \sim \frac{e^{-E_{\theta}(\mathbf{x})}}{Z_{\theta}}$$

Dominated minimum
free-energy
configurations

$$\{\mathbf{x}\}_{\text{eq},\theta} \sim \mathcal{D}$$

$$\left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\theta}} \approx \left\langle \frac{\partial E}{\partial \theta_i} \right\rangle_{p_{\mathcal{D}}} \quad \forall \theta_i$$

Generating new samples

Empirical

Model

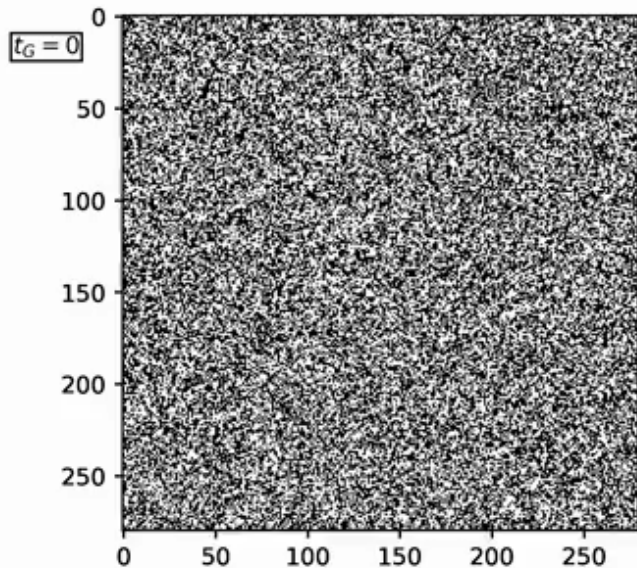
$$p_{\mathcal{D}}(\mathbf{x}) \sim \frac{e^{-E_{\theta}(\mathbf{x})}}{Z_{\theta}}$$

Dominated minimum
free-energy
configurations

$$\{\mathbf{x}\}_{\text{eq},\theta} \sim \mathcal{D}$$



**Markov Chain
Monte Carlo (MCMC)
Langevin dynamics**



Generating new samples

Empirical

Model

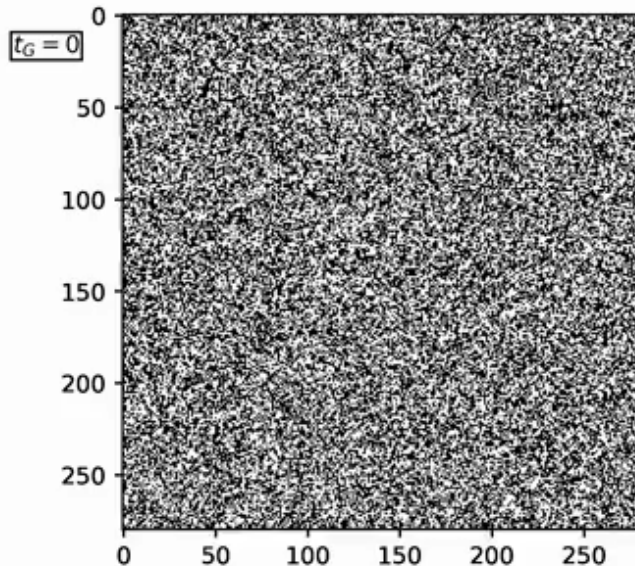
$$p_{\mathcal{D}}(\mathbf{x}) \sim \frac{e^{-E_{\theta}(\mathbf{x})}}{Z_{\theta}}$$

Dominated minimum
free-energy
configurations



**Markov Chain
Monte Carlo (MCMC)
Langevin dynamics**

$$\{\mathbf{x}\}_{\text{eq},\theta} \sim \mathcal{D}$$



$E_{\theta}(\mathbf{x})$ *Effective model
for the data*

If simple, we
can analyze it!

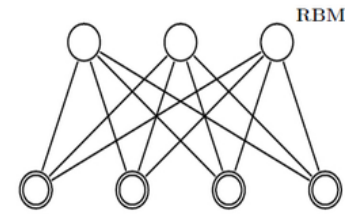
⇒ *Free-energy landscape*

Modeling, interpretability



Why Restricted Boltzmann Machines (RBMs) are good for that?

Why RBMs?



- **Simple enough** to allow some level of analytical treatment (MF)

A Decelle, C Furtlehner - Chinese Physics B, 2021

J Tubiana, R Monasson - Physical review letters, 2017

A Decelle, G Fissore, C Furtlehner - Journal of Statistical Physics, 2018

A Decelle, G Fissore, C Furtlehner

Europhysics Letters, 2017

Biroli, Decelle, Bachtis, Seoane (2024, in prep.)

- Phase diagram

- Learning : sub-sequence of phase transitions

- Approximate methods to compute the free energy (TAP eqs.)

Gabrié, M., Tramel, E. W., & Krzakala, F. NeurIPS (2015)

Tramel, E. W., Gabrié, M., Manoel, A., Caltagirone, F., & Krzakala, F. Physical Review X (2018)

Decelle, A., Rosset, L., & Seoane, B. PRE (2023)

- Can be mapped to a physical interacting system

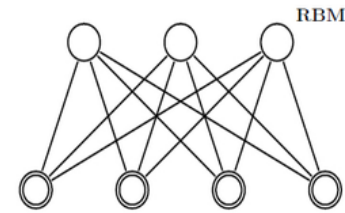
Decelle, Furtlehner, Navas & Seoane, B. SciPost Phys (2024)

- They are **expressive** : they can describe interesting datasets

- They are **frugal** models : fast code and to train

- They are **sample efficient** : perform well with small amounts of data

Why Restricted Boltzmann Machines



- **Simple enough** to allow some level of analytical treatment (MF)

A Decelle, C Furtlehner - Chinese Physics B, 2021

- Phase diagram

**Aurélien Decelle's
lecture tomorrow**

J Tubiana, R Monasson - Physical review letters, 2017

A Decelle, G Fissore, C Furtlehner - Journal of Statistical Physics, 2018

- Learning : sub-sequence of phase transitions

A Decelle, G Fissore, C Furtlehner
Europhysics Letters, 2017

Biroli, Decelle, Bachtis, Seoane (2024, in prep.)

- Approximate methods to compute the free energy (TAP eqs.)

Gabrié, M., Tramel, E. W., & Krzakala, F. NeurIPS (2015)

Tramel, E. W., Gabrié, M., Manoel, A., Caltagirone, F., & Krzakala, F. Physical Review X (2018)

Decelle, A., Rosset, L., & Seoane, B. PRE (2023)

- Can be mapped to a physical interacting system

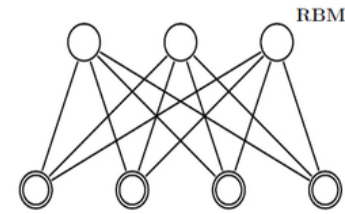
Decelle, Furtlehner, Navas & Seoane, B. SciPost Phys (2024)

- They are **expressive** : they can describe interesting datasets

- They are **frugal** models : fast code and to train

- They are **sample efficient** : perform well with small amounts of data

Why Restricted Boltzmann Machines



- **Simple enough** to allow some level of analytical treatment (MF)

A Decelle, C Furtlehner - Chinese Physics B, 2021

J Tubiana, R Monasson - Physical review letters, 2017

A Decelle, G Fissore, C Furtlehner - Journal of Statistical Physics, 2018

A Decelle, G Fissore, C Furtlehner

Europhysics Letters, 2017

Biroli, Decelle, Bachtis, Seoane (2024, in prep.)

- Phase diagram

- Learning : sub-sequence of phase transitions

- Approximate methods to compute the free energy (TAP eqs.)

Gabrié, M., Tramel, E. W., & Krzakala, F. NeurIPS (2015)

Tramel, E. W., Gabrié, M., Manoel, A., Caltagirone, F., & Krzakala, F. Physical Review X (2018)

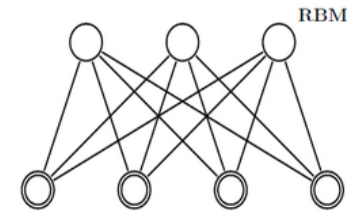
Decelle, A., Rosset, L., & Seoane, B. PRE (2023)

- Can be mapped to a physical interacting system

Decelle, Furtlehner, Navas & Seoane, B. SciPost Phys (2024)

- They are **expressive** : they can describe interesting datasets
- They are **frugal** models : fast code and to train
- They are **sample efficient** : perform well with small amounts of data

Why Restricted Boltzmann Machines



PLOS GENETICS

BROWSE

PUBLISH

ABOUT

SEARCH

advanced search

OPEN ACCESS PEER-REVIEWED

RESEARCH ARTICLE

Creating artificial human genomes using generative neural networks

Burak Yelmen, Aurélien Decelle, Linda Ongaro, Davide Marnetto, Corentin Tallec, Francesco Montinaro, Cyril Furtlehner, Luca Pagani, Flora Jay

109 Save	78 Citation
20,493 View	102 Share

PLOS COMPUTATIONAL BIOLOGY

advanced search

OPEN ACCESS PEER-REVIEWED

RESEARCH ARTICLE

Deep convolutional and conditional neural networks for large-scale genomic data generation

Burak Yelmen, Aurélien Decelle, Leila Lea Boulos, Antoine Szatkownik, Cyril Furtlehner, Guillaume Charpiat, Flora Jay

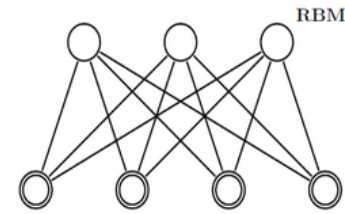
12 Save	2 Citation
1,520 View	15 Share

Version 2 Published: October 30, 2023 • <https://doi.org/10.1371/journal.pcbi.1011584>

Currently the most accurate method to generate artificial human genome

- They are **expressive** : they can describe interesting datasets
- They are **frugal** models : fast code and to train
- They are **sample efficient** : perform well with small amounts of data

Why Restricted Boltzmann Machines



TOOLS AND RESOURCES



Learning protein constitutive motifs from sequence data

Jérôme Tubiana, Simona Cocco, Rémi Monasson*

Laboratory of Physics of the Ecole Normale Supérieure
Research, Paris, France

They are able to capture biologically interpretable features related to function or structure...

Propose mutational paths that can be validated in experiments

PHYSICAL REVIEW LETTERS

Highlights Recent Accepted Collections Authors Referees Search Press About Editorial

Mutational Paths with Sequence-Based Models of Proteins: From Sampling to Mean-Field Characterization

Eugenio Mauri, Simona Cocco, and Rémi Monasson
Phys. Rev. Lett. **130**, 158402 – Published 12 April 2023

Article

References

Citing Articles (3)

Supplemental Material

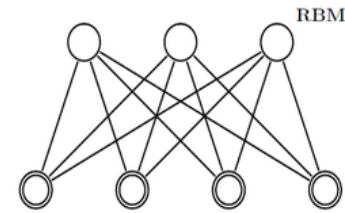
PDF

HTML

Export Citation

- They are **expressive** : they can describe interesting datasets
- They are **frugal** models : fast code and to train
- They are **sample efficient** : perform well with small amounts of data

Why Restricted Boltzmann Machines



- **Simple enough** to allow some level of analytical treatment (MF)
A Decelle, C Furtlehner - Chinese Physics B, 2021
- **Phase diagram**
J Tubiana, R Monasson - Physical review letters, 2017
- **If they are so cool, why are not they used more often?**
A Decelle, G Fissore, C Furtlehner - Europhysics Letters, 2017
- **Learning : sub-sequence of phase transitions**
Biroli, Decelle, Bachtis, Seoane (2024, in prep.)
- **EBMs are very difficult to train properly**
- **Approximate methods to compute the free energy (TAP eqs.)**
Gabrié, M., Tramel, E. W., & Krzakala, F. NeurIPS (2015)
- **Class 2 : Interpretability**
Tramel, E. W., Gabrié, M., Manoel, A., Caltagirone, F., & Krzakala, F. Physical Review X (2018)
- **Can be mapped to a physical interacting system**
Decelle, A., Rosset, L., & Seoane, B. PRE (2023)
- **Class 3: Controlling the training**
Decelle, Furtlehner, Navas & Seoane, B. SciPost Phys (2024)
- They are **expressive** : they can describe interesting datasets
- They are **frugal** models : fast code and to train
- They are **sample efficient** : perform well with small amounts of data