

Département **de Physi**que **École norm**ale

Ecole normale
supérieure





Mean-Field Framework for Unsupervised Learning with Boltzmann Machines

LRI – November 22nd, 2017

Marylou Gabrié Éric Tramel, Andre Manoel, Francesco Caltagirone & Florent Krzakala

Machine Learning, Unsupervised Learning, Density Estimation

Machine learning: design algorithm able to learn from training data

Supervised learning:

learning a mapping between input and output (label)

e.g. classification, regression



Unsupervised learning:

learning without labels, explaining key features of the data

e.g. clustering/community detection, dimensionality reduction, **generative models/density estimation**





Boltzmann Machines fundamental ML model borrowing from Physics

Generative model with pairwise interactions: an Ising model

Set of binary correlated random variables $\mathbf{x} \in \{0,1\}^d$

Boltzmann measure
$$p_{\mathbf{W},\mathbf{b}}(\mathbf{x}) = \frac{e^{-E_{\mathbf{W},\mathbf{b}}(\mathbf{x})}}{Z_{\mathbf{W},\mathbf{b}}}$$

(T. Sejnowski & G. Hinton 1985)



 $\theta = \{\mathbf{W}, \mathbf{b}\}$

Ising energy

$$E_{\mathbf{W},\mathbf{b}}(\mathbf{x}) = -\sum_{i,j} W_{ij} x_i x_j - \sum_i b_i x_i$$
$$E_{\mathbf{W},\mathbf{b}}(\mathbf{x}) = -\mathbf{x}^T \mathbf{W} \mathbf{x} - \mathbf{b}^T \mathbf{x}$$

Boltzmann Machines fundamental ML model borrowing from Physics

Introducing latent variables for representational power

Restricted Boltzmann Machine (RBM) (bipartite)



$$p_{\theta}(\mathbf{v}) = \sum_{\{\mathbf{h}\}} \frac{e^{-E_{\theta}(\mathbf{v},\mathbf{h})}}{Z_{\theta}}$$

Representation learning from hidden state with Deep BM



(R. Salakhutdinov & G. Hinton 2009)

RBM learning requires many approximations

Maximum likelihood learning has no closed form solution

$$\max_{\theta} \prod_{\mathbf{v}^* \in \text{training} \\ \mathbf{v}^* \in \text{training set}} p_{\theta}(\mathbf{v}^*) = \ln p_{\theta}(\mathbf{v}^*) = \ln \left(\sum_{\{\mathbf{h}\}} e^{-E_{\theta}(\mathbf{v}^*, \mathbf{h})} \right) - \ln Z_{\theta}$$

$$data dependent easy log partition hard$$

Iterative gradient ascent training algorithm $W_{i\alpha} \leftarrow W_{i\alpha} + \Delta W_{i\alpha}$

With approximated gradients

$$\Delta W_{i\alpha} \propto \frac{\partial \ell(\theta | \mathbf{v}^*)}{\partial W_{i\alpha}} = \langle v_i^* h_\alpha \rangle_{\mathbf{v}^*} - \langle v_i h_\alpha \rangle$$

$$\bigwedge easy \qquad \bigwedge hard$$

$$p_\theta(\mathbf{h} | \mathbf{v}^*) \quad p_\theta(\mathbf{v}, \mathbf{h})$$

$$\langle \dots \rangle \equiv \mathbb{E}[\dots$$

 h_{α}

 \mathbf{W}

h

Monte Carlo heuristics for RBMs



Contrastive-divergence Monte Carlo with non-thermalized samples

 $\langle v_i h_\alpha \rangle \simeq \frac{1}{K} \sum_{k=1}^K v_k^{(T)} h_k^{(T)}$ with T = O(1)! (G. Hinton 2002)

Successful training but ...

Gibbs sampling unlikely to be accurate for

complicated real valued landscapes

(more sophisticated schemes:

parallel tempering, annealed importance sampling ..)



 $\langle ... \rangle \equiv \mathbb{E}[...]$

A mean-field framework instead ?

Previous works

On RBMs and DBMs

- Tieleman et al. (2009): Naive (or Weiss, or fully factorized) mean-field unsatisfactory
- Salakhutdinov & Hinton (2009): Using naive mean-field for training would be wrong

On fully visible Ising

- Kappen (99): Successful learning with linear response corrections on fully visible BM
- Cocco & Monasson (2012): Cluster expansion including higher orders for inverse Ising problem
- Ricci-Tersenghi (2012): Comparison of the Bethe approximation and other mean-field methods for inverse Ising problem

Can we use extended mean field methods to efficiently train latent variables BMs ?

Thouless-Anderson-Palmer free energy



Intractable distribution and log-partition function = free energy

$$\mathbf{s} \in \{0,1\}^d \qquad E(\mathbf{s}) = -\mathbf{s}^T \mathbf{W} \mathbf{s} - \mathbf{b}^T \mathbf{s} \qquad p_\theta(\mathbf{s}) = \frac{1}{Z_\theta} e^{\mathbf{s}^T \mathbf{W} \mathbf{s} + \mathbf{b}^T \mathbf{s}} \qquad \ln Z_\theta = -F$$

Plefka (82), Georges-Yedidia (99) derivation

Legendre transform $-\beta F(\mathbf{q}) = \ln \sum_{\{\mathbf{s}\}} e^{-\beta (E(\mathbf{s}) + \mathbf{q}^T \mathbf{s})} \Rightarrow \begin{cases} -\beta G(\mathbf{m}) = -\beta \max_{\mathbf{q}} [F(\mathbf{q}) + \mathbf{q}^T \mathbf{m}] \\ \text{Mean-field parameters} = \text{magnetizations} \end{cases}$ Large temperature expansion up to second order $G(\mathbf{m}) \approx G_0 + \beta \left. \frac{\partial G}{\partial \beta} \right|_{\beta=0} + \frac{\beta^2}{2} \left. \frac{\partial^2 G}{\partial \beta^2} \right|_{\beta=0}$ Self consistency $\left. \frac{dG}{dm_i} \right|_{m_i} = 0 \Rightarrow m_i = \sigma \left[\beta b_i + \sum_j \beta W_{ij} m_j - \beta^2 W_{ij}^2 \left(m_i - \frac{1}{2} \right) \left(m_j - m_j^2 \right) \right]_{j \in \mathbb{N}}$

TAP free energy as a function of magnetizations

$$F_{\text{TAP}} = S_{\text{MF}} - \sum_{i} \beta b_{i} m_{i} - \sum_{i,j} \beta W_{ij} m_{i} m_{j} - \sum_{i,j} \beta^{2} \frac{W_{ij}^{2}}{2} (m_{i} - m_{i}^{2}) (m_{j} - m_{j}^{2})$$
NMF

 $\langle ... \rangle \equiv \mathbb{E}[...]$

TAP free energy on the restricted graph

Intractable joint distribution

$$\mathbf{v} \in \{0,1\}^{d^{\mathbf{v}}} \quad \mathbf{h} \in \{0,1\}^{d^{\mathbf{h}}} \quad p_{\theta}(\mathbf{v},\mathbf{h}) = \frac{1}{Z_{\theta}} e^{\mathbf{v}^T \mathbf{W} \mathbf{h} + \mathbf{b}_{\mathbf{v}}^T \mathbf{v} + \mathbf{b}_{\mathbf{h}}^T \mathbf{h}} \qquad \ln Z_{\theta} = -F$$

TAP equations

- magnetizations parameters $\mathbf{m^v} = \langle \mathbf{v} \rangle, \ \mathbf{m^h} = \langle \mathbf{h} \rangle \in [0, 1]$
- small weights expansion of free energy F $(W \leftarrow \beta W)$

$$F_{\text{TAP}} = S_{\text{MF}}(\theta) - \mathbf{b_v}^T \mathbf{m^v} - \mathbf{b_h}^T \mathbf{m^h} - \mathbf{m^v}^T \mathbf{Wm^h} - \sum_{i,\alpha} \frac{W_{i\alpha}^2}{2} (m_i^{\mathbf{v}} - m_i^{\mathbf{v}^2}) (m_{\alpha}^{\mathbf{h}} - m_{\alpha}^{\mathbf{h}^2})$$

$$\mathsf{NMF}$$

$$\mathsf{TAP}$$

self consistency relations for TAP free energy modes

$$m_{i}^{\mathbf{v}} = \sigma \left[b_{\mathbf{v}i} + \sum_{\alpha} W_{i\alpha} m_{\alpha}^{\mathbf{h}} - W_{i\alpha}^{2} \left(m_{i}^{\mathbf{v}} - \frac{1}{2} \right) \left(m_{\alpha}^{\mathbf{h}} - m_{\alpha}^{\mathbf{h}^{2}} \right) \right]$$
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$
TAP

Unsupervised learning with a TAP machine

(M.G, E.Tramel, F. Krzakala NIPS 2015)

We propose to define a TAP machine as:

A bipartite graphical model parametrized by $heta=\{\mathbf{W},\mathbf{b_v},\mathbf{b_h}\}$

Implicitely defining a set of magnetizations $\ \{(m_1^v,m_1^h),(m_2^v,m_2^h),...\}$

Computed as the fixed point of TAP equations

$$m_{i}^{\mathbf{v}} \leftarrow \sigma \left[\mathbf{b}_{\mathbf{v}i} + \sum_{\alpha} W_{i\alpha} m_{\alpha}^{\mathbf{h}} - W_{i\alpha}^{2} \left(m_{i}^{\mathbf{v}} - \frac{1}{2} \right) \left(m_{\alpha}^{\mathbf{h}} - m_{\alpha}^{\mathbf{h}^{2}} \right) \right]$$
$$m_{\alpha}^{\mathbf{h}} \leftarrow \sigma \left[\mathbf{b}_{\mathbf{h}\alpha} + \sum_{i} W_{i\alpha} m_{i}^{\mathbf{v}} - W_{i\alpha}^{2} \left(m_{\alpha}^{\mathbf{h}} - \frac{1}{2} \right) \left(m_{i}^{\mathbf{v}} - m_{i}^{\mathbf{v}^{2}} \right) \right]$$



We define a natural and tractable objective for unsupervised learning: The TAP log-likelihood TAP log-likelihood

$$\ln Z_{\theta} \simeq -F_{\text{TAP}} \quad \Rightarrow \quad \ell_{\text{TAP}}(\theta | \mathbf{v}^*) = \ln \left(\sum_{\{\mathbf{h}\}} e^{-E_{\theta}(\mathbf{v}^*, \mathbf{h})} \right) + F_{\text{TAP}}$$

Gradient ascent of the TAP likelihood: a tractable training scheme for TAP machines



Magnetizations of TAP machines after learning

$$m_{i}^{\mathbf{v}} \leftarrow \sigma \begin{bmatrix} b_{\mathbf{v}i} + \sum_{\alpha} W_{i\alpha} m_{\alpha}^{\mathbf{h}} - W_{i\alpha}^{2} \left(m_{i}^{\mathbf{v}} - \frac{1}{2} \right) \left(m_{\alpha}^{\mathbf{h}} - m_{\alpha}^{\mathbf{h}^{2}} \right) \\ \text{NMF} \end{bmatrix}$$
TAP

TAP fixed points with trained weights and biases for different initializations



Now for a naive-MF machine



Evolution of magnetizations in « phase space »



Isomap Dim. 1

(E.Tramel, M.G., A. Manoel, F. Caltagirone, F. Krzakala 2017)

Evolution of magnetizations in « phase space »



Isomap Dim. 1

(E.Tramel, M.G., A. Manoel, F. Caltagirone, F. Krzakala 2017)

Evolution of magnetizations in « phase space »



(E.Tramel, M.G., A. Manoel, F. Caltagirone, F. Krzakala 2017)

Evolution of magnetizations in « phase space »



Isomap Dim. 1

(E.Tramel, M.G., A. Manoel, F. Caltagirone, F. Krzakala 2017)

Evolution of magnetizations in « phase space »



Isomap Dim. 1

(E.Tramel, M.G., A. Manoel, F. Caltagirone, F. Krzakala 2017)

Mean-Field Framework for Unsupervised Learning

1 – TAP formalism to define TAP machines

- Deterministic model inspired by Boltzmann Machines
- Endowed with a tractable training algorithm
- Acheiving similar quality of training than CD
- Providing new ways of monitoring learning (TAP likelihood, magnetizations)

2 - An application of TAP machines to Bayesian inference

Learning complex structured priors for compressed sensing (CS)

Bayesian inference problem



Machine learning construction of priors for Bayesian inference

- Learn an RBM from data samples to estimate underlying P₀
- But for an RBM P₀(x) is intractable
- Resort instead to a MF latent variable model

Prior learning for compressed sensing



Learn to locate the non-zero entries of x with a TAP machine for the reconstruction ?

Learning the support of MNIST digits for CS

y = Fx + w

Priors tested

Binarized

MNIST

r) E

99

- (a) Gauss-Bernouli i.i.d.:
- (b) Gauss-Bernoulli non-i.i.d.:
- NMF latent graph. model: (c)
- (d) TAP latent graph. Model:

(E.Tramel, A. Drémeau, F. Krzakala J. Stat. Mech 2016)

$$P_0(x_i) = (1 - \rho)\delta(x_i) + \rho \mathcal{N}(x_i; \mu, \sigma^2)$$

$$P_0(x_i) = (1 - \rho_i)\delta(x_i) + \rho_i \mathcal{N}(x_i; \mu, \sigma^2)$$
$$m_i^{\mathbf{v}} = \sigma \left[\theta_{v_i} + \sum W_{i,v_i} m_i^{\mathbf{h}}\right]$$

$$m_{i}^{\mathbf{v}} = \sigma \begin{bmatrix} \theta_{\mathbf{v}i} + \sum_{\alpha} W_{i\alpha} m_{\alpha}^{\mathbf{h}} \\ \theta_{\mathbf{v}i} + \sum_{\alpha} W_{i\alpha} m_{\alpha}^{\mathbf{h}} - W_{i\alpha}^{2} \left(m_{i}^{\mathbf{v}} - \frac{1}{2} \right) \left(m_{\alpha}^{\mathbf{h}} - m_{\alpha}^{\mathbf{h}^{2}} \right) \end{bmatrix}$$

Factor graph representation of posterior $P(\mathbf{x}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{x})P_0(\mathbf{x})$ Fully factorized priors MF latent variable priors





An efficient way to incorporate ML to CS ?

Priors tested

- (a) Gauss-Bernouli i.i.d.
- (b) Gauss-Bernoulli non-i.i.d.

- (c) NMF inference on graph. model
- (d) TAP inference on graph. model

Increasing number of measurements M

Increasing number of measurements M



Mean-Field Framework for Unsupervised Learning

- **1 TAP formalism to define TAP machines**
- 2 An application of TAP machines to Bayesian inference
 - Structured prior can be incorporated thanks to MF/TAP inference on MF/TAP machines.
 - Reconstruction is possible from a lot fewer measurements.
- 3 Generalization to arbitrary variables and deep networks

Generalized Restricted Boltzmann Machines

(E.Tramel, M.G., A. Manoel, F. Caltagirone, F. Krzakala arxiv.org/abs/1702.032602017)

Arbitrary continuous priors for real valued units

$$\mathbf{v} \in \mathbb{R}^{d^{\mathbf{v}}} \qquad p_{\theta}(\mathbf{v}, \mathbf{h}) = \frac{1}{Z_{\theta}} e^{\mathbf{v}^{T} \mathbf{W} \mathbf{h}} \prod_{i} P_{\theta_{i}^{\mathbf{v}}}(v_{i}) \prod_{\alpha} P_{\theta_{\alpha}^{\mathbf{h}}}(h_{\alpha}) \qquad \mathbf{W}$$

Log – likelihood

$$\ell(\theta|\mathbf{v}^*) = \ln p_{\theta}(\mathbf{v}^*) = \sum_{\alpha} \ln \int dh_{\alpha} P_{\theta_{\alpha}^{\mathbf{h}}}(h_{\alpha}) e^{\sum_{i} W_{i\alpha} v_{i}^* h_{\alpha}} + \sum_{i} \ln P_{\theta_{i}^{\mathbf{v}}}(v_{i}^*) - \ln Z_{\theta}$$

TAP approximation using « mean » and « variance » variational parameters

This again recovers Approximate Message Passing algorithm

$$\mathbf{m}^{\mathbf{v}} = \langle \mathbf{v} \rangle, \ \mathbf{c}^{\mathbf{v}} = \langle \mathbf{v}^2 \rangle - \langle \mathbf{v} \rangle^2$$

$$\mathbf{m}^{\mathbf{h}} = \langle \mathbf{h} \rangle, \ \mathbf{c}^{\mathbf{h}} = \langle \mathbf{h}^2 \rangle - \langle \mathbf{h} \rangle^2$$

$$F_{\text{TAP}}(\mathbf{m}^{\mathbf{v}}, \mathbf{c}^{\mathbf{v}}, \mathbf{m}^{\mathbf{h}}, \mathbf{c}^{\mathbf{h}})$$

Training continuous TAP machines

CBCL Faces



Distributions

Binary hidden

$$P_0(h_{\alpha}) = \frac{e^{\theta_{\alpha}^{\mathbf{h}} h_{\alpha}}}{1 + e^{\theta_{\alpha}^{\mathbf{h}}}}$$

Truncated Gaussian visibles

$$P_0(v_i) = \mathcal{N}(v_i; \theta_i^{\mathbf{v}}, \sigma_i^2) \mathbb{1}_{0 \le v_i \le 1}$$

Log-likelihood

 $\ell_{\mathrm{TAP}}(\mathbf{W}, \theta^{\mathbf{h}}_{\alpha}, \theta^{\mathbf{v}}_{i}, \sigma^{2}_{i} | \mathbf{v})$

Learned weights



Marginal means

$$\mathbf{m^v} = \langle \mathbf{v}
angle$$









Isomap Dim. 1







Isomap Dim. 1



Isomap Dim. 1



2

Isomap Dim.

Isomap Dim. 1

Prior learning and inference with real values



Binary model for the support only (BRBM)

$$\mathbf{v} \in \{0, 1\}_{d^{\mathbf{h}}}^{d^{\mathbf{v}}} p_{\theta}(\mathbf{v}, \mathbf{h}) = \frac{1}{Z_{\theta}} e^{\mathbf{v}^{T} \mathbf{W} \mathbf{h}} \prod_{i} e^{\theta_{i}^{\mathbf{v}^{T}} v_{i}} \prod_{\alpha} e^{\theta^{\mathbf{h}^{T}} h_{\alpha}}$$
$$\mathbf{h} \in \{0, 1\}^{d^{\mathbf{v}}} p_{\theta}(\mathbf{v}, \mathbf{h}) = \frac{1}{Z_{\theta}} e^{\mathbf{v}^{T} \mathbf{W} \mathbf{h}} \prod_{i} e^{\theta_{i}^{\mathbf{v}^{T}} v_{i}} \prod_{\alpha} e^{\theta^{\mathbf{h}^{T}} h_{\alpha}}$$

Generalized model on real valued digits (GRBM)

$$\mathbf{h} \in \{0, 1\}^{d^{\mathbf{h}}} \quad p_{\theta}(\mathbf{v}, \mathbf{h}) = \frac{1}{Z_{\theta}} e^{\mathbf{v}^{T} \mathbf{W} \mathbf{h}} \prod_{i} P_{\theta_{i}^{\mathbf{v}}}(v_{i}) \prod_{\alpha} e^{\theta_{\alpha}^{\mathbf{h}^{T}} h_{\alpha}} \\ \mathbf{v} \in \mathbb{R}^{d^{\mathbf{v}}} \quad P_{\theta_{i}^{\mathbf{v}}}(v_{i}) = (1 - \rho_{i}) \delta(v_{i})^{i} + \rho_{i} \mathcal{N}(v_{i}; \overset{\alpha}{\mu}_{i}, \sigma_{i}^{2})$$



Deep Boltzmann machine learning

Data dependent terms becomes intractable because of direct interactions between hidden units

log-likelihood

 $log \text{ partition hard} \qquad data \text{ dependent easy}$ $\ell(\theta|\mathbf{x}) = \ln P(\mathbf{x};\theta) = -\ln \mathcal{Z}[\theta] + \sum_{i} \ln P_{i}(x_{i};\theta_{i})$ $+ \ln \int \left(\prod_{l} \prod_{j} dh_{j}^{(l)}\right) e^{\sum_{i,j} x_{i} W_{ij} h_{j}^{(1)} + \sum_{l} \sum_{i,j} h_{i}^{(l)} W_{ij}^{(l)} h_{j}^{(l+1)}} \prod_{l} \prod_{j} P_{j}^{h(l)}(h_{j}^{(l)};\theta_{j}^{(l)})$ data dependent hard !

h₃

 h_2

• gradients w.r.t. trained parameters

$$\Delta W^{(l)}{}_{i\alpha} \propto \frac{\partial \ell(\theta | \mathbf{x}^*)}{\partial W_{i\alpha}} = \langle h_i^{(l-1)} h_{\alpha}^{(l)} \rangle_{\mathbf{x}^*} - \langle h_i^{(l-1)} h_{\alpha}^{(l)} \rangle$$

$$\uparrow \mathsf{hard}$$

$$p_{\theta}(\mathbf{h} | \mathbf{v}^*) \qquad p_{\theta}(\mathbf{v}, \mathbf{h})$$

Deep Boltzmann machine learning

$$\Delta W^{(l)}{}_{i\alpha} \propto \frac{\partial \ell(\theta | \mathbf{x}^*)}{\partial W_{i\alpha}} = \langle h_i^{(l-1)} h_\alpha^{(l)} \rangle_{\mathbf{x}^*} - \langle h_i^{(l-1)} h_\alpha^{(l)} \rangle$$

Usual training algorithm:

- NMF approximation of data dependent terms
- MCMC approximation of second term

The mean field framework suggests instead the consistent scheme

- TAP approximation of data dependent term
- TAP approximation of second term





Conclusion & Perspectives

- TAP equations allow for new unsupervised learning with latent variables
 - New way of monitoring an unsupervised training
 - New way of using ML for Bayesian inference
- Mean field framework deals with arbitrary variables and architecture

PERSPECTIVES

- Bayesian inference problems with very meaning full priors ?
- Further analytical analysis of the TAP machine ? Spectral analysis ?

THANK YOU !