Unsupervised learning of features from data: a statistical physics approach

R. Monasson Laboratory of Physics CNRS & Ecole Normale Supérieure, Paris

Model-guided data science, Como, September 2019

Plan of the lectures

1. Bayesian Inference and dimensional reduction: phase transition in principal component analysis

2. Representations: auto-encoders, Restricted Boltzmann Machines & sparse feature learning

3. Restricted Boltzmann Machines: connections with graphical models, phase transitions & applications

Restricted Boltzmann Machines

• **Graphical model** constituted by two sets of random variables that are coupled together.

$$P(v,h) = \frac{1}{Z} \exp\left[-E(v,h)\right]$$

$$E(v,h) = -\sum_{i} g_{i}v_{i} + \sum_{\mu} U_{\mu}(h_{\mu}) - \sum_{i,\mu} w_{i\mu}v_{i}h_{\mu}$$



Smolensky 1986

Restricted Boltzmann Machines: Learning

• **Graphical model** constituted by two sets of random variables that are coupled together.

$$P(v,h) = \frac{1}{Z} \exp\left[-E(v,h)\right]$$

$$E(v,h) = -\sum_{i} g_{i}v_{i} + \sum_{\mu} U_{\mu}(h_{\mu}) - \sum_{i,\mu} w_{i\mu}v_{i}h_{\mu}$$

• RBM learns a **probability distribution** over the **visible layer**.



$$P(v) = \int \prod_{\mu} dh_{\mu} P(v, \{h_{\mu}\}) \equiv \frac{1}{Z_{eff}} \exp\left[-E_{eff}(v)\right]$$

• We also have P(h | v) extract (distribution of) representations from data P(v | h) generate (distribution of) data given a representation

RBM are generative models, trained through unsupervised learning

Hidden-Unit potential

• Compute hidden units inputs
$$I_{\mu}^{H} = \sum_{i} W_{i\mu} V_{i}$$

• Sample hidden units $P(h_{\mu} | I_{\mu}^{H}) \propto \exp\left[-U_{\mu}(h_{\mu}) + h_{\mu} I_{\mu}^{H}\right]$

Most likely value of h given input I^{H} ?

$$\frac{d}{dh}U(h) = I^{H} \Longrightarrow h = \Phi(I^{H})$$



 $\mathbf{v}^{(t)}$

Bernoulli or ReLU empirically known to give better results than linear hidden units ...

MNIST synthetic digits: Linear vs. ReLU RBMs



MNIST 60,000 image of digits of size 28x28 <u>Linear RBM</u> M = 400 LogLikelihood : -0.15 bits/ pixel <u>ReLU RBM</u> M = 400 LogLikelihood : -0.11 bits/pixel

(digits are less noisy, more accurate, and MCMC are mixing faster)

RBM vs. BM: quadratic hidden potentials

Quadratic potential for hidden units:

$$U(h) = \frac{h^2}{2}$$

 Energy of visible and hidden configurations:

2. Joint distribution:

$$E(v,h) = -\sum_{i} g_{i}v_{i} + \sum_{\mu} U_{\mu}(h_{\mu}) - \sum_{i,\mu} w_{i\mu}v_{i}h_{\mu}$$

$$P(v,h) = \frac{1}{Z} \exp\left[-E(v,h)\right]$$

3. Marginal distribution over visible configurations

$$P(v) = \int \prod_{\mu} dh_{\mu} P(v, \{h_{\mu}\}) = \frac{1}{Z_{eff}} \exp\left[-E_{eff}(v)\right]$$

RBM vs. BM: quadratic hidden potentials

$$E_{eff}(v) = -\sum_{i} g_{i}v_{i} + \frac{1}{2}\sum_{\mu} \left(\sum_{i} w_{i\mu}v_{i}\right)^{2} = -\sum_{i} g_{i}v_{i} + \frac{1}{2}\sum_{i,j} \left(\sum_{\mu} w_{i\mu}w_{j\mu}\right)v_{i}v_{j}$$

$$\uparrow$$

$$J_{ij}$$

Boltzmann machine, i.e. Ising model with interaction matrix of rank-M ! Also called Hopfield model

Non quadratic potentials generate multi-body interactions between v_i ...



 g_i

Talk by B. Bravi next Friday

Hopfield's model (1982)

Proc. Natl. Acad. Sci. USA Vol. 79, pp. 2554–2558, April 1982 Biophysics

Neural networks and physical systems with emergent collective computational abilities

(associative memory/parallel processing/categorization/content-addressable memory/fail-soft devices)

J. J. HOPFIELD

Division of Chemistry and Biology, California Institute of Technology, Pasadena, California 91125; and Bell Laboratories, Murray Hill, New Jersey 07974

Contributed by John J. Hopfield, January 15, 1982

Computational properties of use to biological or-ABSTRACT ganisms or to the construction of computers can emerge as collective properties of systems having a large number of simple equivalent components (or neurons). The physical meaning of content-addressable memory is described by an appropriate phase space flow of the state of a system. A model of such a system is given, based on aspects of neurobiology but readily adapted to integrated circuits. The collective properties of this model produce a content-addressable memory which correctly yields an entire memory from any subpart of sufficient size. The algorithm for the time evolution of the state of the system is based on asynchronous parallel processing. Additional emergent collective properties include some capacity for generalization, familiarity recognition, categorization, error correction, and time sequence retention. The collective properties are only weakly sensitive to details of the modeling or the failure of individual devices.

- Autoassociative memory
- Simple dynamics of components (no clock)
- Generalization, error correction, time sequence storage, ...
- Robustness to failure of individual components

The model

set of activity configurations (patterns) to be 'stored':

Index of pattern = 1, ..., P

Index of neuron = 1, ..., N



synaptic interactions: $J_{ij} = \frac{1}{N} \sum_{\mu} w_i^{\mu} w_j^{\mu}$

updating rule: $v_i(t+1) = sign\left(\sum_i J_{ij}v_j(t) - \theta_i\right)$ $(\theta_i = -g_i)$

Q: Is the final state close to one of the patterns?

A: Yes, if number P of patterns small enough ...



Are the minima of E close to the patterns?

Reminder on the mean field-Ising model

Probability of a configuration of spins:

$$P_{K}(v_{1}, v_{2}, ..., v_{N}) = \frac{1}{Z(K)} \exp\left(\frac{K}{2N} \sum_{i \neq j} v_{i} v_{j}\right)$$
$$m(K) = \sum P(v_{1}, v_{2}, ..., v_{N}) \times \frac{1}{2N} \sum P(v_{1}, v_{N}) \times \frac{1}{2N} \sum$$

Order parameter:

$$m(K) = \sum_{\{v_1, v_2, \dots, v_N = \pm 1\}} P(v_1, v_2, \dots, v_N) \times \frac{1}{N} \sum_i v_i$$

characterizes the degree of order of a typical configuration:



From the Mattis model...



Same phase transition!

... to the Hopfield model



Same phase transition as Ising model if p finite & N tends to infinity:

i.e., for large K, one overlap m^{μ} is positive, the others vanish (other scenarios with multiple positive overlaps are exponentially unlikely...)

The Hopfield model in the double $p, N \rightarrow \infty$ limit



Gaussian RBM: The Phases (at low temperature)



Prototype vs. compositional regimes in nonlinear RBM







16 hidden units (prototypes)

100 hidden units (sparse features)

Fischer & Igel. Training Restricted Boltzmann Machines: An Introduction, 2014.

Prototype regime



Strongly reminiscent of ferromagnetic regime in Hopfield model :



- One hidden unit is strongly activated and the others are weakly activated;
- Limited diversity of visible-layer configurations given hidden-unit activation
- Number of high-probability configurations is linear in nb. of hidden units



Compositional regime







$0.06 \\ 0.04 \\ 0.02 \\ 0.00 \\ 0.02 \\ 0.00 \\ L \\ 0.00 \\ L \\ 0.00 \\ S \\ 0.00 \\ 0.$

Compositional regime

Each generated digit image is composed by superposing about L≈20 elementary strokes, while S≈250 units are silent. Different combinations of strokes produce different variants of the same digit or different digits



A subset of the features learnt for a ReLU RBM. M = 400.

Each image represent a weight vector

$$\left\{ W_{i}^{}\right\} _{\mu}^{}$$



-0.50.00.5



- After learning of data, weights are sparse
- And few, remaining weights are large!



Conditional averages Black: v=0, White: v=1, Grey scale

• No interference thanks to large ReLU threshold ...

Random-Weight RBM ensemble



Statistical Mechanics of Random-Weight RBMs

• Replica theory computation is performed to estimate the average free energy *F* in the zero-temperature limit :



The three representational regimes of RBM

"Ferromagnetic" regime



- One hidden unit very active
- Corresponding weights define local prototype
- Unable to extract invariances

More hidden units

"Glassy" regime



- All hidden units active
- Visible configurations are complex mixtures
- Not Interpretable

Increasing sparsity

Non quadratic hidden-unit potentials



Compositional regime

- Multiple hidden units very active
- Corresponding weights define features composing visible configurations
- Extract invariances; interpretable

Transition paths





Example of Transition Path

(diffusive motion)

Proteins: structure & function



Sequence fully defines 3D structure and function of the corresponding protein: Genotype-phenotype mapping?

Proteins: structure & function



- Sequence fully defines 3D structure and function of the corresponding protein: Genotype-phenotype mapping?
- Mapping is complex: Mutations may largely affect stability, activity, specificity, ... of protein (directed evolution for protein design) No additivity a priori: epistasis

Proteins: Evolutionary Families



- Many realizations of sequences of proteins with similar function & structure
- Characterization of underlying distribution??
- May help to Infer structure, function; Predict effect of mutations; Design new proteins; ...

Multi-sequence alignments

Q9U3M7_CAEEL/420-758	YQ		-YPRYVTE-I	TVRHLLEHTA	GGWDNLQSDA
062481_CAEEL/42-406	YLPEF	РА	KKFKNEDVKI	TMRQLLSHSA	GIRHYATEKK
Q9I1G0_PSEAE/16-386	WLPAF	RP	RLADGREARI	TPRQLLSHSA	GLGYRFLEAD
ESTB_BURGA/22-391	WLPEF	RP	RLADGSEPLV	TIHHLLTHTS	GLGYW-LLEG
LOVD_ASPTE/17-406	RLLPDLSAMP	VLEGFDDAGN	ARLRERRGKI	TLRHLLTHTS	GLSYVFLHPL
Q9WXD6_BRELN/11-387	YAPGLADVQV	I-EGFDVDGS	PILRAPASEP	TTKQLLLHTA	GFGYDFFNEK
Q9L8D3_SORCE/35-419	WLPELANRKV	LARIDGPI	DETVPAERPI	TVRDLMTFTM	GFGISFDASS
Q9A2D7_CAUCR/97-472	YLPEFADM	KVSDGQ	GGVRPAARPI	LVRDILRHTA	GFSYGWGEGP
Q9A800_CAUCR/77-466	HIPEFANLRV	A-KGVDETGQ	PILTPVSRSP	TMRELMTHTA	GFAYGLANDP
Q9A7Z9_CAUCR/50-429	FIPEFANLRV	L-KGVNADGS	FDTVPAERPP	TMRELMSHSA	GFAYGLTPDN
088012_STRC0/23-405	YLPAFAEPRV	YVGGTGEN	VVTRPATGPV	RVRHLMTHTA	GLTFGFYRTH
Q9A3L3_CAUCR/99-483	FIPAWRDIGV	FQAGVAGA	FQTTRTKEPM	RTIDLLRHTS	GLTYGFQQRT
P72041_MYCTU/40-404	YLP		SYTSHGKHRT	TIRHVLTHSA	GVPFPTGPRP
Q9X5Q0_STRLA/21-361	YWP		QFARHGKGDV	TVRHVLQHRA	GVPVGRGIVR
Q18384_CAEEL/50-416	YWP		EFGQNGKQDI	TIEMVLSHTA	GLPYFPGVKF
Q9XU43_CAEEL/52-422	YWP		EYGRYGKNAT	TIEDVLSHKA	GLPYL-SEDI

Sequence number

Amino-acid content

[- = gap, same format for all sequences]

High-dimensional representations of protein sequences



- RBM extract high-D representations of (common inputs to) sequences
- Representations are useful (to design « good » sequences) ...
- ... and, hopefully, biologically meaningful (structure, function, history)
 - ➔ Practical implementation of genotype-to-phenotype relation

Applications of RBM to protein sequence data

WW domain (PFAM PF00397)



Tubiana, Cocco, R.M., eLife 2019

Example 1: the WW domain

- Binding domain involved in eukaryotic signalling proteins
- N=30-40 amino-acids (very small)
- Folds into 3-stranded antiparallel beta strands



- Recognizes Proline-rich Linear Motifs with 4 types of ligand specificities
 - Type I: PPXY, Y = Tyrosine (aromatic) , X = any residue
 - Type II: PPLP, L = Leucine
 - Type III: PR rich peptide, R = arginine
 - Type IV: p(S/T)P, phosphorylated serine/threonine

Example 2: HSP70 chaperone protein

- N>600 amino-acids
- Multidomain.
 - Nucleotide Binding Domain (NBD)
 - Substrate Binding Domain (SBD)
 - LID Domain
 - Linker

Function:

- Traps substrate proteins between the LID and the SBD
- LID/SBD cavity is either open or close

Roles:

- Assist protein folding
- Transport proteins for degradation



ATP bound conformation (open) PDB: 2kho ADP bound conformation (closed) PDB: 4jne

Hidden-unit potentials



Generalization of Rectified Linear Units:

- Positive slopes at origin favors h=0 (as L1 regularization)
- Negative slopes favors bimodal distribution for h
- Confining potential at large h
- Four parameters to be learned from data for each hidden unit

Weights reflect structural features

WW domain









Weights reflect structural features HSP70



- Collective mode located on the unstructured C- terminal tail
- Known to be crucial for interactions with co-chaperones
- Either charged hydrophilic or tiny hydrophobic residues
- Analogous to IDP

Interdomain weights control allostery



Nucleotide Binding Domain (NBD) LID Domain Substrate Binding Domain (SBD) Linker

Interdomain features control allostery



Loop motifs control ATP/ADP regulation



RBM features reflect function













WW Domain: Weights Determine Motif Recognition Specificity



Sequence

space

The problem:

Find probability distribution from very few samples

Mixture of local models :

Each hidden unit sees and codes for a patch in sequence space



Sequence space



All or almost all hidden units active at any position in sequence space

Non interpretable representations ...



Sequence space

Decomposition into constitutive features:

Each hidden codes for an invariant feature; sequences are obtained by combinatorial composition of features



The three representational regimes of RBM

"Ferromagnetic" regime



- One hidden unit very active
- Corresponding weights define local prototype
- Unable to extract invariances

More hidden units

"Glassy" regime



- All hidden units active
- Visible configurations are complex mixtures
- Not Interpretable

Increasing sparsity

Non quadratic hidden-unit potentials



Compositional regime

- Multiple hidden units very active
- Corresponding weights define features composing visible configurations
- Extract invariances; interpretable

Driving RBM to the compositional phase



Driving RBM to the compositional phase



WW Domain: Weights Determine Motif Recognition Specificity



Artificial Sequence Generation with RBM



Artificial Sequences

Artificial Sequence Generation with RBM



Artificial Sequences

Type II-like binding pocket + Short loop



Type II-like binding pocket + Short loop



Artificial Sequence Generation with RBM



Type II/III/IV-like binding pocket + Short loop → Type II/III



Type II/III/IV-like binding pocket + Short loop → Type IV



Perspectives

- Data-driven models are getting increasingly important; need for controlled approaches to infer and interpret models.
- If you have questions about lectures or notes, please contact me! (remi.monasson at ens.fr)
- Fascinating issues from statistical, computational, conceptual points of views. Lots of things to do in future!!
- Postdoc positions available from January 2020 ...