

Marc Mézard

Ecole normale supérieure PSL University

Leiden, Colloquium Ehrenfestii, November 14, 2018

What is inference?

Brain, deduction, reasoning, learning...

Psychology : inference = production of new internal representations on the basis of previously held representations

eg : - perceptions \rightarrow expectations \rightarrow action

- conceptual thinking

Fusion: posterior parietal

regions

Segregation: primary visual & auditory cortices

Causal inference: anterior parietal regions

Not necessarily conscious (while reasoning is)

see eg Mercier Sperber 2011

What is inference?

Statistics

Infer a hidden rule, or hidden variables, from data. Restricted sense : find parameters of a probability distribution

Urn with 10.000 balls. Draw 100, find 70 white balls and 30 black Best guess for the composition of the urn? How reliable? Probability that it has 6000 white- 4000 black?

If only black and white balls , with fraction x of white, probability to pick-up 70 white balls is $\binom{100}{70}x^{70}(1-x)^{30}$

Log likelihood of x: $L(x) = 70 \log x + 30 \log(1 - x)$ Maximum at $x^* = .7$ Probability of .6 : $e^{L(.6) - L(.7)}$

Bayesian inference

Unknown parametersxPriorP(x)MeasurementsyLikelihoodP(y|x)

Posterior $P(x|y) = \frac{P(y|x)P(x)}{P(y)}$

Bayesian inference

Unknown parameters	x	Prior	P(x)
Measurements	y	Likelihood	P(y x)

Posterior $P(x|y) = \frac{P(y|x)P(x)}{P(y)}$

Bayesian inference

Unknown parametersxPriorP(x)MeasurementsyLikelihoodP(y|x)

Posterior

 $P(\mathbf{x}|y) = \frac{P(y|\mathbf{x})P(\mathbf{x})}{P(y)}$



What is inference?

Artificial intelligence, machine learning



Find a machine that reads handwritten digits...

... inferring its parameters from examples



MNIST database : 70,000 images of digits, segmented, 28 × 28 pixels each, greyscale. Known output (supervised learning) ₆

What is inference?

Artificial intelligence, machine learning



« Neural network » : artificial neurons



$$y = f(w_0 + w_1x_1 + w_2x_2 + w_3x_3)$$

Formal neural network





What is inference?

Artificial intelligence, machine learning



Machine with hundreds of thousands of parameters, trained on very large data base: infer the parameters from data (supervised learning)

What is inference?

Information theory, communication, signal processing

Information transfer : error correction by the use of redundancy





















Statistical inference

Challenge = rules with **many hidden parameters**. eg : machine learning with large machine and big data, decoding in commonication,...

$$x = (x_1, \dots, x_N) \quad N \gg 1$$

Many measurements $y = (y_1, \dots, y_M)$ $M \gg 1$

Measure of the amount of data $\alpha = M/N$

Algorithms

Prediction on the quality of inference, on the performance of the algorithms, on the type of situations where they can be applied

Bayesian inference with many unknown and many measurements

Unknown parameters
$$x = (x_1, \dots, x_N)$$
Prior $P^0(x)$ Measurements $y = (y_1, \dots, y_M)$ $P(y|x)$

Bayesian inference

$$P(x|y) \propto P(y|x)P^0(x)$$

Often (but not necessarily): Independent measurements

$$P(y|x) = \prod_{\mu} P_{\mu}(y_{\mu}|x)$$

Factorized prior
$$P^{0}(x) = \prod_{i} P_{i}^{0}(x_{i})$$

Posterior $P(x) = \frac{1}{Z(y)} \left(\prod_{i} P_{i}^{0}(x_{i})\right) \exp\left[-\sum_{\mu} E_{\mu}(x, y_{\mu})\right]$

 $E_{\mu}(x, y_{\mu}) = -\log P_{\mu}(y_{\mu}|x)$

Bayesian inference with many unknown

and many measurements

$$P(x) = \frac{1}{Z(y)} \left(\prod_{i} P_i^0(x_i) \right) \exp \left[-\sum_{\mu} E_{\mu}(x, y_{\mu}) \right]$$

$$E_{\mu}(x, y_{\mu}) = -\log P_{\mu}(y_{\mu}|x)$$

Statistical mechanics.

« Spin glass »

\bigstarDiscrete or continuous variables x_i

♦Interactions through $e^{-E_{\mu}(x,y_{\mu})}$ can be

•pairwise :
$$E_{\mu} = J_{\mu} x_i(\mu) x_j(\mu)$$

multibody

Disordered system, ensemble
Thermodynamic limit, phase transitions



 $s_i = \pm 1$

• Disordered magnetic systems

e.g.: CuMn





$$P(s_1,\ldots,s_N) = \frac{1}{Z}e^{-E/T}$$

Spin glasses

• Disordered magnetic systems e.g.

e.g.: CuMn



Each spin 'sees' a different local field

 Many atoms, microscopic interactions are known, "disordered systems" e.g.: CuMn



Each spin 'sees' a different local field
 Low temperature: frustration



 Many atoms, microscopic interactions are known, "disordered systems"
 e.g.: CuMn



Each spin 'sees' a different local field
Low temperature: frustration
Spins freeze in random directions
Difficult to find min. of E





Many quasi-ground states unrelated by symmetries, many metastable states

Slow dynamics, aging

Spin glass

Each spin 'sees' a different local field
Low temperature: frustration
Spins freeze in random directions
Difficult to find min. of E





Many quasi-ground states unrelated by symmetries, many metastable states

Slow dynamics, aging

Spin glass

Each spin 'sees' a different local field
Low temperature: frustration
Spins freeze in random directions
Difficult to find min. of E

Useless, but thousands of papers...



Inference, spinglass and crystal: tomography of binary mixtures



Inference, spinglass and crystal: L measurements for each angle tomography of binary mixtures α Synchrotron light





If the size of domains is \gg pixel: possible to reconstruct with $\ll L^2$ measurements

 $\xi \gg 1$



If the size of domains is \gg pixel: possible to reconstruct with $\ll L^2$ measurements

 $\xi \gg 1$

Tomography of binary mixtures

Compressed

sensing

This picture, digitalized on 1000 × 1000 grid, can be reconstructed fom measurements with 16 angles

> Gouillart et al., Inverse problems 2013

If the size of domains is \gg pixel: possible to reconstruct with $\ll L^2$ measurements



 $\mu \qquad y_{\mu} = \sum s_i$ $i \in \partial \mu$

Prior knowledge on $\{s_i\}$: neighboring pixels more likely to be equal

Studied with mean-field



 $\boldsymbol{\mu} \qquad y_{\mu} = \sum_{i \in \partial \mu} s_i$

Prior knowledge on $\{s_i\}$: neighboring pixels more likely to be equal

$$P(S) = \prod_{ij \in \text{grid}} e^{Js_i s_j} \prod_{\mu} \delta\left(y_{\mu}, \sum_{i \in \partial \mu} s_i\right)$$

Studied with mean-field

prior

measurement

$$P(S) = \prod_{ij \in \text{grid}} e^{Js_i s_j} \prod_{\mu} \delta\left(y_{\mu}, \sum_{i \in \partial \mu} s_i\right)$$

If enough measurements: The most probable S (the ground state) gives the perfect composition of the sample.

« Crystal » : much more probable

But in some cases « crystal hunting » may be computationally very hard !





« Crystal » : much more probable

But in some cases « crystal hunting » may be computationally very hard !



Inference with many unknowns : « crystal hunting » with mean-field based algorithms Historical development of mean field equations :

- In homogeneous ferromagnets:
 - Weiss (infinite range, 1907)
 - Bethe Peierls (finite connectivity, 1935)
- In glassy systems:
 - Thouless Anderson Palmer 1977,
 - M. Parisi Virasoro 1986 (infinite range)
 - M. Parisi 2001 (finite connectivity)

- As an algorithm:
- Gallager 1963
- Pearl 1986
 - M. Parisi Zecchina 2002
 - Kabashima 2003, 2008
 - Donoho Bayati Montanari 2009
 - Rangan 2010
 - Krzakala M. Zdeborova 2012 ..
BP = Bethe-Peierls = Belief Propagation



 $P(x_1, \cdots, x_5) = \psi_a(x_1, x_2, x_4)\psi_b(x_2, x_3)\cdots$



First type of messages:

Probability of x_1 in the absence of a:

 $m_{1 \rightarrow a}(x_1)$



Second type of messages:

Probability of x_1 when it is connected only to c:

$$m_{c \to 1}(x_1)$$









Propagate messages along the edges, update messages at vertices, using elementary local probabilistic rules



Propagate messages along the edges, update messages at vertices, using elementary local probabilistic rules

 $m_{3 \to g}(x_3)$

Closed set of equations: two messages "propagate" on each edge of the factor graph. When is BP exact?

$$m_{1 \to c}(x_1) = Cm_{d \to 1}(x_1)m_{e \to 1}(x_1)m_{f \to 1}(x_1)$$
$$m_{c \to 2}(x_2) = \sum_{x_1, x_3} \psi_c(x_1, x_2, x_3)m_{1 \to c}(x_1)m_{3 \to c}(x_3)$$

Fluctuations are handled correctly, but beware of correlations

- Exact in one dimension (transfer matrix
 - = dynamic programming)
- Exact on a tree (uncorrelated b.c)
- Exact on locally tree-like graphs (Erdös Renyi etc.) if correlations decay fast enough (single pure state) and uncorrelated disorder
- Exact in infinite range problems if correlations decay fast enough (single pure state) and uncorrelated disorder



Two important developments

1) The special case of infinite-range models

2) What happens in a glass phase, when there are many pure states, and therefore many solutions ?

Infinite range models



eg TAP equations for spin glasses:

$$m_{i \to \mu}(x_i) = \prod_{\nu \neq \mu} m_{\nu \to i}(x_i)$$
$$M(x_i) = \prod_{\nu \neq \mu} m_{\nu \to i}(x_i)$$

$$M_i(x_i) = \prod_{\nu} m_{\nu \to i}(x_i)$$

Small difference, treated perturbatively Mean-field equations can be written only in terms of site pdfs: $M_i(x_i)$. TAP, AMP....

$$t+1 \qquad t \qquad t-1 \qquad t-1 \qquad t-1 \qquad t-1 \qquad t-1$$
$$H_i = \sum_k J_{ki} \tanh(\beta H_k) - \beta \tanh(\beta H_i) \sum_k J_{ki}^2 [1 - \tanh^2(\beta H_k)]$$

Two important developments

1) The special case of infinite-range models

2) What happens in a glass phase, when there are many pure states, and therefore many solutions ?

2) What happens in a glass phase, when there are many pure states, and therefore many solutions ?

BP equations

$$m_{i \to \mu}(x_i) = \prod_{\nu \neq \mu} m_{\nu \to i}(x_i)$$

Correct if, in absence of the i-j interaction, the correlations between k and ℓ can be neglected.



Loop length $O(\log N)$

2) What happens in a glass phase, when there are many pure states, and therefore many solutions ?

BP equations

$$m_{i \to \mu}(x_i) = \prod_{\nu \neq \mu} m_{\nu \to i}(x_i)$$

Correct if, in absence of the i-j interaction, the correlations between k and ℓ can be neglected.



Energy

$$m_{i \to \mu}^{\alpha}(x_i) = \prod_{\nu \neq \mu} m_{\nu \to i}^{\alpha}(x_i)$$

Configurations

Glassy phase: many states, many solutions of BP

Loop length $O(\log N)$

2) What happens in a glass phase, when there are many pure states, and therefore many solutions ?

BP equations

$$m_{i \to \mu}(x_i) = \prod_{\nu \neq \mu} m_{\nu \to i}(x_i)$$

Correct if, in absence of the i-j interaction, the correlations between k and ℓ can be neglected.

$$\bigvee_{\nu \to \mu} \mathcal{M}_{i \to \mu}(x_i) = \prod_{\nu \neq \mu} m^{\mathcal{A}}_{\nu \to i}(x_i)$$

Statistics of $m_{i \to \mu}^{\alpha}(x_i)$

over the many states α

$$P_{i \to \mu}(m)$$

related to

$$P_{\nu \to i}(m)$$

Survey propagation M Parisi Zecchina 2002

Configurations

Glassy phase: many states, many solutions of BP

Power of message passing algorithms

Approximate solution of very hard, and very large constraint satisfaction problems, ...FAST! (typically linear time)

- BP: Best decoders for LDPC error correcting codes
- SP: Best solver of random satisfiability problems
- BP: Best algorithm for learning patterns in neural networks (e.g. binary perceptron)
- Data clustering, graph coloring, Steiner trees, etc...
- Fully connected networks : TAP (=AMP). Compressed sensing, linear estimation, etc.

Local, simple update equations:
Each message is updated using
information from incoming
messages on the same node.
Distributed, solves hard global pb

Power of message passing algorithms

Approximate solution of very hard, and very large constraint satisfaction problems, ...FAST! (typically linear time)

- BP: Best decoders for LDPC error correcting codes
- SP: Best solver of random satisfiability problems
- BP: Best algorithm for learning patterns in neural networks (e.g. binary perceptron)
- Data clustering, graph coloring, Steiner trees, etc...
- Fully connected networks : TAP (=AMP). Compressed sensing, linear estimation, etc.





An example of fully connected model: Generalized Linear Regression

Unknowns: x_i $i = 1, \ldots, N$

Linear combinations:
$$z_{\mu} = \sum_{i} F_{\mu i} x_{i}$$
 $\mu = 1, ..., M$
Outputs y_{μ} generated from $P_{out}(y_{\mu}|z_{\mu})$

Prior factorized

$$\int P(x_i)$$

Bayes

$$P(x|y) = \frac{1}{Z(y)} \prod_{i} P(x_i) \prod_{\mu} P_{\text{out}}(y_{\mu}|\sum_{i} F_{\mu i}x_i)$$

Examples: tomography, linear regression, perceptron learning, compressed sensing...



Linear regression: Individual μ : expression of disease y_{μ}

Value of factor i for individual $\mu : F_{\mu i}$

Find the best weights of factors x_i Minimize mean square error with regularization

$$\frac{1}{2}\sum_{\mu}(y_{\mu}-\sum_{i}F_{\mu i}x_{i})^{2}+\sum_{i}||x_{i}||$$

$$P_{out} = e^{-(y_{\mu} - z_{\mu})^2/(2\Delta)}$$

$$P(x_i) = e^{-||x_i||/\Delta}$$

$$P(x|y) = \frac{1}{Z(y)} \prod_{i} P(x_i) \prod_{\mu} P_{\text{out}}(y_{\mu}|\sum_{i} F_{\mu i}x_i)$$

Compressed sensing

Unknown variables x_i

Linear measurements
$$y_{\mu} = \sum_{i} F_{\mu i} x_{i} + \eta_{\mu}$$

Compressed sensing regime : $M < N$
 $P_{out} = e^{-(y_{\mu} - z_{\mu})^{2}/(2\Delta)}$

sparse prior (in appropriate basis)

$$P(x_i) = (1 - \rho)\delta(x_i) + \rho\phi(x_i)$$





$$P(x|y) = \frac{1}{Z(y)} \prod_{i} P(x_i) \prod_{\mu} P_{\text{out}}(y_{\mu}|\sum_{i} F_{\mu i}x_i)$$

 $F_{\mu i}$: iid, known

Spin glass with multispin interactions, infinite range: write mean field equations.

TAP equations written in terms

of
$$a_i = \langle x_i \rangle$$

 $c_i = \langle x_i^2 \rangle - \langle x_i \rangle^2$

Iteration — algorithm : AMP Statistical study — phase diagram

 $F_{\mu i}$ x_i $P(x_i)$

 $P_{\rm out}(y_{\mu}|z_{\mu})$

M 1989, OW 96, K 2003, K 2008 , DMM 2009, R2011, KMSSZ 2012

Benchmark: noiseless limit of compressed sensing with iid measurements

System of linear measurements



Random F : «random projections» (incoherent with signal) Pb: Find x when M < N and x is sparse

Phase diagram

«Thermodynamic limit»

 $N \gg 1$ variables $R = \rho N$ non-zero variables $M = \alpha N$ equations

• Solvable by enumeration when $\alpha > \rho$ but $O(e^N)$

• ℓ_1 norm approach Find a *N* - component vector *x* such that the *M* equations y = Fx are satisfied and $||x||_1$ is minimal • AMP = Bayesian approach Planted: $\phi_T(x)$

$$P(\mathbf{x}) = \prod_{i=1}^{N} [(1-\rho)\delta(x_i) + \rho\phi(x_i)] \prod_{\mu=1}^{P} \delta\left(y_{\mu} - \sum_{i} F_{\mu i} x_i\right)$$

Performance of AMP with Gauss-Bernoulli prior: phase diagram



Analysis of random instances : phase transitions

N (real) variables, M measurements (linear functions)

Analysis of random instances : phase transitions

Reconstruction of signal using BP. Fixed $\,^{\rho}$, decrease $\,^{\alpha}$





Dynamical phase transition. Ubiquitous in statistical inference. Conjecture « All local algorithms freeze »... How universal?

Step 3: design the measurement matrix in order to get around the glass transition

Getting around the glass trap: design the matrix F so that one nucleates the naive state (crystal nucleation idea,

...borrowed from error correcting codes : « spatial coupling »)

Felström-Zigangirov, Kudekar Richardson Urbanke, Hassani Macris Urbanke,

«Seeded BP»

Nucleation and seeding



Nucleation and seeding





 $F_{\mu i}$ = independent random Gaussian variables, zero mean and variance $J_{b(\mu)b(i)}/N$



F

whole system!

- L = 8 $\alpha_1 > \alpha_{BP}$ $N_i = N/L$ $M_i = \alpha_i N / L$
 - $\alpha_{j} = \alpha' < \alpha_{BP} \qquad j \ge 2$ $\alpha = \frac{1}{L} \left(\alpha_{1} + (L-1)\alpha' \right)$

S



Performance of the probabilistic approach + message passing + parameter learning+ seeding matrix

$$Z = \int \prod_{j=1}^{N} \mathrm{d}x_j \prod_{i=1}^{N} \left[(1-\rho)\delta(x_i) + \rho\phi(x_i) \right] \prod_{\mu=1}^{M} \delta\left(y_\mu - \sum_{i=1}^{N} F_{\mu i} x_i \right)$$



Simulations
Analytic approaches (replicas and cavity)

$$\rightarrow \alpha_c = \rho_0$$

Reaches the ultimate information-theoretic threshold

Proof: Donoho Javanmard Montanari

Performance of AMP with Gauss-Bernoulli prior: phase diagram





Phase transitions are crucial in large inference problems Hard-Impossible = absolute limit (Shannon-like) Easy- Hard = limit for large class of algorithms (local)

The spin glass cornucopia

A very sophisticated and powerful corpus of conceptual and methodological approaches has been developed (replicas, cavity, TAP,...) mostly in the years 1975-2000, and has found applications in many different fields of

information theory and computer science

Portrait of Ottavio Strada, Tintoretto, Venice 1567 Rijk's Museum Amsterdam


Back to error correction Efficient codes : parity checks (LDPC codes)

Add redundancy, with structure allowing to decode

 $x_i \in \{0, 1\}$



 $a: x_1 + x_4 + x_5 + x_7 = 0 \pmod{2}$

- $b: x_2 + x_4 + x_6 + x_7 = 0 \pmod{2}$
- $c: x_3 + x_5 + x_6 + x_7 = 0 \pmod{2}$

 2^4 codewords among 2^7 words



$$P(x_1, \dots, x_N | y_1, \dots, y_N) = \frac{1}{Z} \prod_i \psi_i(x_i | y_i) \prod_a \mathbb{I}\left(\sum_{i \in \partial a} x_i = 0 \pmod{2}\right)$$

Spin glass problem with multispin interactions, discontinuous glass transition (1 step RSB)



 $P(x_1, \dots, x_N | y_1, \dots, y_N) = \frac{1}{Z} \prod_i \psi_i(x_i | y_i) \prod_a \mathbb{I}\left(\sum_{i \in \partial a} x_i = 0 \pmod{2}\right)$ received

Spin glass problem with multispin interactions, discontinuous glass transition (1 step RSB)

Error decoding: « crystal hunting » inference problem

 $P(x_1, \dots, x_N | y_1, \dots, y_N) = \frac{1}{Z} \prod_i \psi_i(x_i | y_i) \prod_a \mathbb{I}\left(\sum_{i \in \partial a} x_i = 0 \pmod{2}\right)$ received

A priori knowledge of the channel

Spin glass problem with multispin interactions, discontinuous glass transition (1 step RSB)

Error decoding: « crystal hunting »
inference problem

$$A priori knowledge of$$

$$P(x_1, \dots, x_N | y_1, \dots, y_N) = \frac{1}{Z} \prod_i \psi_i(x_i | y_i) \prod_a \mathbb{I} \left(\sum_{i \in \partial a} x_i = 0 \pmod{2} \right)$$

A priori knowledge of the channel Parity check constraints

Spin glass problem with multispin interactions, discontinuous glass transition (1 step RSB)

Error decoding: inference problem



$$P(x_1, \dots, x_N | y_1, \dots, y_N) = \frac{1}{Z} \prod_i \psi_i(x_i | y_i) \prod_a \mathbb{I}\left(\sum_{i \in \partial a} x_i = 0 \pmod{2}\right)$$

One possible decoding algorithm: use belief-propagation mean-field equations relating the local fields

Solve them iteratively (Gallager)

Phase Transitions in Error correcting codes

Shannon 1948 (random code ensemble)

Typical structured code **ensemble** (e.g. LDPC), with optimal decoding

Typical structured code **ensemble** , with fast BPbased decoding algorithm



Phase Transitions in Error correcting codes

Shannon 1948 (random code ensemble)

Typical structured code **ensemble** (e.g. LDPC), with optimal decoding

Typical structured code **ensemble** , with fast BPbased decoding algorithm



Phase Transitions in Error correcting codes

Shannon 1948 (random code ensemble)

Typical structured code **ensemble** (e.g. LDPC), with optimal decoding

Typical structured code **ensemble** , with fast BPbased decoding algorithm

















0









The spin glass cornucopia

A very sophisticated and powerful corpus of conceptual and methodological approaches has been developed (replicas, cavity, TAP,...) mostly in the years 1975-2000, and has found applications in many different fields of

information theory and computer science

Portrait of Ottavio Strada, Tintoretto, Venice 1567 Rijk's Museum Amsterdam

