

From neural networks to graphical models: a brief introduction

Toshiyuki Tanaka ¹

Graduate School of Informatics, Kyoto University,

36-1 Yoshida-honmachi, Sakyo-ku, Kyoto-shi, Kyoto 606-8501, Japan

1 Neural networks and graphical models

One of the most important and interesting aspects of the neural networks is that they exhibit collective information-processing capabilities by connecting processing elements (“neurons”), each of which performs a very simple and limited information processing. The central issue of the neural network research is thus to explore how and why such a system with simple elements exhibits complex information processing functionalities. Being looked at from another direction, by studying the problem, we are asking about the global characterization of the information processing performed by the network, on the basis of the local characterization of information processing (i.e., that performed by each neuron), as well as how they are interconnected.

Essentially the same problem also arises in the study of graphical models (Bayesian networks). The main objective of this section is to give a brief review of the graphical models, with some emphasis on their relation to neural networks. In one important class of graphical models, a system with an N -dimensional random vector \mathbf{x} is characterized by a set of local conditional probabilities of the form: $p(x_i|\mathbf{x}_{\setminus i})$, where $\mathbf{x}_{\setminus i}$ denotes the set of all the elements of \mathbf{x} except x_i . If, for all i , x_i is independent of x_j , $j > i$, then the set of the local conditional probabilities defines a global probability distribution

$$p(\mathbf{x}) = \prod_{i=1}^N p(x_i|x_1, \dots, x_{i-1}). \quad (1)$$

If, on the other hand, the above independence condition does not hold,

$$\rho(\mathbf{x}) = \prod_{i=1}^N p(x_i|\mathbf{x}_{\setminus i}) \quad (2)$$

does not necessarily define a global probability distribution, but with the proper normalization coefficient $Z_0 = \sum_{\mathbf{x}} \rho(\mathbf{x})$, one can define a global probability distribution by $p(\mathbf{x}) = \rho(\mathbf{x})/Z_0$. It should be noted that such a model can be regarded as a generalization of the conventional neural networks with stochastic neurons: In fact, if one assumes that the “state” x_i of the “neuron” i depends on $\mathbf{x}_{\setminus i}$ only through their

¹E-mail: tt@i.kyoto-u.ac.jp

weighted sum $\mathbf{w}_i \cdot \mathbf{x}_{\setminus i}$, with \mathbf{w}_i the weight vector of the neuron i , then the system defined by (2) reduces to a stochastic neural network.

A further generalization is possible by introducing clique functions $p(\mathbf{x}_r)$, where a clique r is a subset of the index set $\{1, \dots, N\}$. We let

$$\rho(\mathbf{x}) = \prod_{r \in C} p(\mathbf{x}_r), \quad (3)$$

where the product is taken over a set of cliques C . Again, $\rho(\mathbf{x})$ does not necessarily define a global probability distribution, but it can be normalized to give a global probability distribution, as $p(\mathbf{x}) = \rho(\mathbf{x})/Z_0$, where

$$Z_0 = \sum_{\mathbf{x}} \rho(\mathbf{x}). \quad (4)$$

As a demonstrative example, we would like to note that (3) defines a Boltzmann machine [1] when the elements of \mathbf{x} are binary and when one considers “two-body” cliques $r = (i, j) \subset \{1, \dots, N\}$. See also [2] for an early example of graphical models that appeared in the neural network literature.

2 Expectations, likelihoods and statistical mechanics

With a model such as those described as (2) or (3), it is often the case that one has to deal with the following problems. Let $\mathbf{x} = (\mathbf{x}_V, \mathbf{x}_H)$, where \mathbf{x}_V and \mathbf{x}_H denote “visible” and “hidden” variables, respectively. Assume that one has an observation \mathbf{a} for the visible variables \mathbf{x}_V . Then, for example, one is interested in evaluating

$$\sum_{\mathbf{x}_H} x_i p(\mathbf{x}_H | \mathbf{x}_V = \mathbf{a}) = \frac{\sum_{\mathbf{x}_H} x_i \rho(\mathbf{x}_H, \mathbf{x}_V = \mathbf{a})}{Z}, \quad (5)$$

that is, expectation of $x_i \in \mathbf{x}_H$ conditioned on the observation $\mathbf{x}_V = \mathbf{a}$. Here we let

$$Z = \sum_{\mathbf{x}_H} \rho(\mathbf{x}_H, \mathbf{x}_V = \mathbf{a}). \quad (6)$$

Such quantities are important when one wants to perform inference on the basis of the observation \mathbf{a} and the probability model (2) or (3) at hand. Or, one would like to evaluate

$$p(\mathbf{x}_V = \mathbf{a}) = \sum_{\mathbf{x}_H} p(\mathbf{x}_H, \mathbf{x}_V = \mathbf{a}) = \frac{Z}{Z_0}, \quad (7)$$

that is, the likelihood of the observation, which measures “goodness” of the model (2) or (3) in explaining the observation \mathbf{a} , and therefore is useful in learning as well as in model selection. Evaluation of conditional expectations and likelihoods is often computationally hard. It is indeed the case when the elements of \mathbf{x}_H and \mathbf{x} are discrete, in which case one generally requires computation that scales exponentially in the dimensions of \mathbf{x}_H and of \mathbf{x} to evaluate Z and Z_0 , respectively.

One can notice formal similarity between the formulation just introduced and those found in statistical mechanics. The quantities Z_0 and Z can be regarded as partition functions of properly defined physical

systems. This is the basic observation behind the idea of applying statistical-mechanics notions and tools to evaluate expectations and likelihoods. If no clique function becomes 0, then one can let

$$\rho(\mathbf{x}) = \exp\left[\sum_{r \in C} c_r(\mathbf{x}_r)\right], \quad (8)$$

which is exactly the Gibbs-Boltzmann distribution with “Hamiltonian” $-\sum_{r \in C} c_r(\mathbf{x}_r)$.

3 Advanced mean field methods

3.1 Naive mean field approximation

In the following, we assume for simplicity that the ranges of the elements of \mathbf{x} are subsets of \mathbb{R} , the set of real values. In statistical mechanics, various approximation frameworks to evaluate expectations and likelihoods have been studied. The simplest one is the naive mean field approximation, in which one considers a “test” distribution of the form

$$p_0(\mathbf{x}; \boldsymbol{\theta}) = \exp\left[\sum_i \theta_i x_i - \psi_0\right], \quad (9)$$

where $\boldsymbol{\theta} = (\theta_1, \dots, \theta_N)$ is the parameter of p_0 , and where ψ_0 is the logarithm of the normalization coefficient of p_0 , which is identified with the free energy in statistical mechanics. One then tries to find the parameter value that gives the “best” approximation to the true distribution $p(\mathbf{x})$. In order to obtain the best approximation we minimize the Kullback-Leibler divergence

$$D(p_0 \| p) = \sum_{\mathbf{x}} p_0(\mathbf{x}; \boldsymbol{\theta}) \log \frac{p_0(\mathbf{x}; \boldsymbol{\theta})}{p(\mathbf{x})} = \boldsymbol{\theta} \cdot \boldsymbol{\eta} - \psi_0 - \sum_{r \in C} \langle c_r(\mathbf{x}_r) \rangle_0 + \psi, \quad (10)$$

where ψ is the free energy of the true distribution p , where $\langle \cdot \rangle_0$ denotes expectation with p_0 , and where $\boldsymbol{\eta} = \langle \mathbf{x} \rangle_0$. The stationarity condition gives the so-called mean field equations

$$\theta_i = \sum_{r \in C} \frac{\partial \langle c_r(\mathbf{x}_r) \rangle_0}{\partial \eta_i}. \quad (11)$$

With a solution of the mean field equations, one can evaluate expectations with p_0 , which hopefully will give a good approximation to the expectations with p . The approximate free energy ψ_0 can also be evaluated from the solution, which is expected to give a good approximation to the true free energy ψ .

An information-geometrical interpretation of the naive mean field approximation is given in [3].

3.2 Belief propagation

In the naive mean field approximation we have considered a single test distribution. Alternatively, we can consider many of them, anticipating that it improves approximation. Specifically, we let

$$p_r(\mathbf{x}; \boldsymbol{\zeta}_r) = \exp\left[\sum_i \zeta_{ri} x_i + c_r(\mathbf{x}_r) - \psi_r\right], \quad r \in C. \quad (12)$$

We then consider the following algorithm:

1. Initialize $\{\boldsymbol{\xi}_r = (\xi_{r1}, \dots, \xi_{rN}), r \in C\}$.

2. Iterate the following:

- (a) For $r \in C$, let $\zeta_r = \sum_{s \in C \setminus r} \xi_s$.
- (b) Compute the naive mean field approximation of $p_r(\mathbf{x}; \zeta_r)$, to obtain $p_0(\mathbf{x}; \theta_r)$.
- (c) Calculate the “extrinsic information” $\xi_r = \theta_r - \zeta_r$.

The algorithm defined above is the belief propagation [4]. Belief propagation gives the correct answer when the graphical representation [4], induced by the cliques, of the true distribution $p(\mathbf{x})$ is a tree. The algorithm is still applicable to cases where the graphical representation has loops, but in such cases we lose theoretical guarantee that it converges to the true solution. Furthermore, it may not converge at all. Nevertheless, the algorithm has been studied extensively, because of its importance in coding theory [5, 6], as well as its relation with the Bethe approximation in statistical mechanics [7, 8].

If the algorithm arrives at an equilibrium, θ_r becomes independent of r . Let the fixed-point value of θ_r be θ^* . Then the following properties hold.

- θ^* is the minimizer of the Kullback-Leibler divergence

$$D(p_0 \| p_r) = \sum_{\mathbf{x}} p_0(\mathbf{x}; \theta) \log \frac{p_0(\mathbf{x}; \theta)}{p_r(\mathbf{x}; \zeta_r^*)} \quad (13)$$

for fixed ζ_r^* .

- The true distribution $p(\mathbf{x})$ and the test distributions $p_0(\mathbf{x}; \theta^*)$, $\{p_r(\mathbf{x}; \zeta_r^*), r \in C\}$ are in a log-linear relation, that is,

$$\log p(\mathbf{x}) - \log p_0(\mathbf{x}; \theta^*) = \sum_{r \in C} [\log p_r(\mathbf{x}; \zeta_r^*) - \log p_0(\mathbf{x}; \theta^*)] + \text{const.} \quad (14)$$

These properties have an elegant information-geometrical interpretation, which is discussed in [9].

In the context of statistical mechanics, “two-body” cliques are usually considered, but one can use more general cliques, for example, spanning trees of the model [10].

References

- [1] D. H. Ackley, G. E. Hinton, and T. J. Sejnowski, “A learning algorithm for Boltzmann machines,” *Cognitive Science*, vol. 9, pp. 147–169, 1985.
- [2] G. E. Hinton and T. J. Sejnowski, “Optimal perceptual inference,” *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 448–453, 1983.
- [3] T. Tanaka, “Information geometry of mean field theory,” *IEICE Trans. Fundamentals*, vol. E79-A, no. 5, pp. 709–715, 1996.
- [4] J. Pearl, *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*, Morgan-Kaufmann, 1988.
- [5] R. J. McEliece, D. J. C. MacKay, and J.-F. Cheng, “Turbo decoding as an instance of Pearl’s ‘belief propagation’ algorithm,” *IEEE J. Select. Areas Commun.*, vol. 16, no. 2, pp. 140–152, 1998.

- [6] D. J. C. MacKay, “Good error-correcting codes based on very sparse matrices,” *IEEE Trans. Info. Theory*, vol. 45, no. 2, pp. 399–431, 1999; Errata, *ibid.*, vol. 47, no. 5, p. 2101, 2001.
- [7] Y. Kabashima and D. Saad, “Belief propagation *vs.* TAP for decoding corrupted messages,” *Europhys. Lett.*, vol. 44, no. 5, pp. 668–674, 1998.
- [8] J. S. Yedidia, W. T. Freeman, and Y. Weiss, “Constructing free-energy approximations and generalized belief propagation algorithms,” *IEEE Trans. Info. Theory*, vol. 51, no. 7, pp. 2282–2312, 2005.
- [9] S. Ikeda, T. Tanaka, and S. Amari, “Stochastic reasoning, free energy, and information geometry,” *Neural Computation*, vol. 16, no. 9, pp. 1779–1810, 2004.
- [10] M. J. Wainwright, T. S. Jaakkola, and A. S. Willsky, “Tree-based reparameterization framework for analysis of sum-product and related algorithms,” *IEEE Trans. Info. Theory*, vol. 49, no. 5, pp. 1120–1146, 2003.