



Perspectives on Probabilistic Graphical Models

DONG LIU

Doctoral Thesis in Electrical Engineering
Stockholm, Sweden 2020

Devision of Information Science and Engineering
TRITA-EE XXXX KTH, School of Electrical Engineering and and Computer Science
ISSN SE-100 44 Stockholm
ISBN SWEDEN

Akademisk avhandling som med tillstånd av Kungl Tekniska högskolan framlägges
till offentlig granskning för avläggande av teknologie doktorsexamen i Elektroteknik
onsdag den 13 september 2017 klockan 13.15 i F3, Lindstedtsvägen 26, Stockholm.

© 2020 Dong Liu, unless otherwise noted.

Tryck: Universitetsservice US AB

To my beloved

Abstract

Sammanfattning

Acknowledgements

Contents

Abstract	v
Sammanfattning	vii
Acknowledgements	ix
Contents	ix
Acronyms and Notations	xi
1 Introduction	1
1.1 Motivations	1
1.2 Scope and Thesis Outline	4
2 Background	9
2.1 Graphical Models	9
2.2 Divergence	13
2.3 Inference Tasks	14
2.4 Variational inference	15
2.5 Learning principles	17

I	Inference	21
3	An alternative view of belief propagation	23
3.1	α belief propagation	23
3.2	Convergence study	23
3.3	Experimental results	23
3.4	Summary	23
4	Region-based Energy Neural Network Model	25
4.1	Region-based graph and energy	25
4.2	RENN model for Approximate Inference	25
4.3	RENN model for markov random field training	25
4.4	Experimental results	25
4.5	Summary	25
II	Learning	27
5	Learning with inference	29
5.1	learning Undirected graphical models/ MRF	29
5.2	Amortized/Neural Variational Learning and Inference of partial observed MRF	29
5.3	Notation	30
5.4	Model and Problem Definition	30
5.5	A lower bound of the marginal likelihood	31
5.6	Experiment	31
6	Powering the expectation maximization method by neural networks	33
6.1	Normalizing flow	34
6.2	expectation maximization of neural network based mixture models	34
6.3	An alternative construction method	34
6.4	Experiments	34
6.5	Summary	34
7	Powering Hidden Markov Model by Neural Network based Generative Models	35
7.1	Hidden Markov Model	35
7.2	GenHMM	35
7.3	Application to phone recognition	35
7.4	Application to sepsis detection in preterm infants	35
7.5	Summary	35
8	An implicit probabilistic generative model	37
8.1	Modeling data without explicit probabilistic distribution	37

8.2	Employing EOT for modeling	37
8.3	Experimental results	37
8.4	Summary	37
 IIIEpilogue		39
9 Conclusion and Discussions		41
Bibliography		43

Acronyms and Notations

Notations

X	random variable
x	realization of the random variable X
\mathcal{X}	alphabet of the random variable X
X_i^k	random sequence (X_i, \dots, X_k)
x_i^k	realization of the random sequence X_i^k
\mathcal{X}_i^k	alphabet of the random sequence X_i^k
X^k	random sequence (X_1, \dots, X_k)
x^k	realization of the random sequence X^k
\mathcal{X}^k	alphabet of the random sequence X^k
$X_i^{k \setminus n}$	random sequence $(X_i, \dots, X_{n-1}, X_{n+1}, \dots, X_k)$
$x_i^{k \setminus n}$	realization of the random sequence $X_i^{k \setminus n}$
$\mathcal{X}_i^{k \setminus n}$	alphabet of the random sequence $X_i^{k \setminus n}$
$X^{k \setminus n}$	random sequence $(X_1, \dots, X_{n-1}, X_{n+1}, \dots, X_k)$
$x^{k \setminus n}$	realization of the random sequence $X^{k \setminus n}$
$\mathcal{X}^{k \setminus n}$	alphabet of the random sequence $X^{k \setminus n}$
$ \cdot $	set cardinality
f_X	p.d.f. of the continuous random variable X
p_X	p.m.f. of the discrete random variable X
$\mathcal{N}(\mu, \sigma^2)$	normal distribution with mean μ and variance σ^2

$D(\cdot \cdot)$	Kullback-Leibler divergence
$D_\tau(\cdot \cdot)$	τ -th order Rényi divergence
$C(\cdot, \cdot)$	Chernoff information
$E[\cdot]$	expectation
$\partial\cdot$	boundary of a closed set
$\hat{\partial}\cdot$	upper boundary of a two-dimensional closed set
$\check{\partial}\cdot$	lower boundary of a two-dimensional closed set
$\log(\cdot)$	natural logarithm

Chapter 1

Introduction

1.1 Motivations

Most tasks conducted by a person or an automated system requires a fundamental ability of *reasoning*, which is always about reaching a conclusion based on available information. At times, a conclusion is not enough and it is also required to know how reliable the conclusion is. Take the coronavirus (COVID-19) that started from Wuhan in China at the end of 2019 as example, a doctor needs to check the information about a person to reason if the person is infected by the coronavirus. The relevant information includes symptoms such as fever, cough, breathing difficulties and probably kidney failure in severe cases. Even after the doctor has reached a conclusion of positive or negative infection of coronavirus for the person, a natural question is why and how *confident* the diagnose is made.

Instead of randomly guessing, reasoning is to answer queries with preserving uncertainty, by making best of available information. Two fundamental problems are inevitable before a rational reasoning can be conducted.

- How should we specify the relationship between a conclusion and the available information? In the coronavirus example, the counterpart question to answer is how the doctor should relate coronavirus infection with the symptoms. This step is called *modeling* which represents a reasoning problem abstractly by specifying the relationship between known information and unknown parts, in preparation of answer query on it.
- With the modeled representation, how a conclusion should be made? This process of reaching a answer to the query from an abstracted representation is called *inference*. Assume the doctor knows that the candidate has fever and cough, but not any breathing problem, how likely does the candidate have be infected by coronavirus?

In the modeling and inference process, it is not likely that the very beginning assumptions are perfectly right about the truth and can be used for all instances of

the same type of queries. On the contrary, we usually begin with simple (sometimes naive) assumptions of representation a problem, and come back to revise it later when more information or evidence is available, which is aligned with our learning process of new knowledge. In fact, this assumption-and-revision procedure can be more compact. Instead of fixing the model representation at the beginning when one might not be sure the correctness of the assumptions about the model, one can assume a set of models with the assumption that the 'correct' model is in this set. A typical strategy is to leave some freedom to the configuration of the model at the beginning. Each instantiating the configuration generates a model representation. By using available observations or information, the model is adjusted to be able to be best compatible with the observations. This procedure adds the following fundamental problem in reasoning.

- Instead of having a fixed model at the first step, a set of models (or hypothesis) is given. We then need to choose one model that best describes the available observations or information. This phase of choosing a model is called *learning*.

Afterwards, inference can be conducted on the learned model to answer queries.

With all the discussed problems above, modeling, inference and learning, our purpose is to carry out reasoning with being aware of how confident a conclusion or answer is. These problems can be treated nicely with probabilistic models via Bayesian logic. Probabilistic model is built on the fundamental calculus of probability theory that is natural to accommodate the *uncertainty*, which is desired in reasoning. In addition, the probabilistic models offers rich space to modeling problems, where inference can be carried on either exactly or approximately. **More importantly, the modeling or modeling learning part is not necessarily coupled with inference algorithms.** This proper separation allows freedom that a certain family of general inference algorithms can be applies to a broad class probabilistic models. It offers the freedom of trying different model representation of a class without the need of replacing inference algorithm.

Example 1. *Consider the coronavirus infection problem. Using probabilistic model, we are able to model the problem in a rigid way. Additionally, we can make query more formally in probabilistic model framework. Assume each symptom among fever, cough and breathing difficulty can take value from {True,False}. Also the coronavirus infection is either true or false. One exemplified query can be*

$$P(\text{Infection} = \text{True} | \text{Fever} = \text{True}, \text{Cough} = \text{False}, \text{BreathingDifficulty} = \text{True}),$$

which is asking how likely the patient is infected by coronavirus if symptoms of both fever and breathing difficulty are observed but no sight of cough symptom.

Given the fact that probabilistic theory offers a rigid foundation to model and study the problems, which is used to answer query that we concern, it soon becomes intractable when dozens or hundreds of relevant attributes are joint considered.

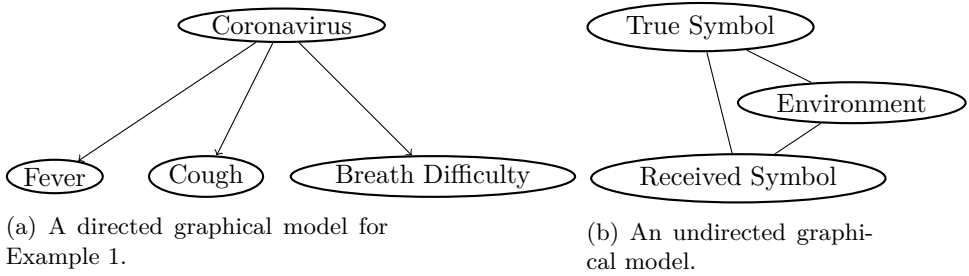


Figure 1.1: Different perspectives on probabilistic graphic models. 1.1a A toy Bayesian network. 1.1b A toy Markov random field.

This can be exemplified by giving finer levels of each symptom in coronavirus infection, e.g. symptom fever is represented by the actual body temperature in integer instead of true-or-false binary state on the one hand. On the other hand, there could be more directly and indirectly relevant symptoms such as muscle pain and congestion. Together with the symptoms, the recent travel itinerary is also related. Additionally, season flu could also similarly bring up some symptoms listed above.

Probabilistic graphical model offers a general framework to encode random variable dependency of a complex probabilistic distribution into a structured graph, which is a powerful tool to compactly model relevant attributes and facts of a complex problem or a system. As show in Figure 1.1a that represents the problem of Example 1 into a directed graphical model (also called Bayesian network or generative model), the nodes (or vertexes) correspond to the variables that represents symptoms and infection state, whereas the edges between nodes correspond how one variable may influence others. In certain scenarios, it is natural to use directed graphical models (generative models) to represent that a observable variable is dependent on a latent variable that generates the observations. For instance, a noisy location sensor of a car keeps measuring the car’s true location and reading noisy locations.

In contract to the directed graphical model, there are more scenarios that the interaction between related random variables is not directional and an undirected edge is used, which leads to the undirected graphical model (or Markov random field, Section 2.1) representation. Undirected graphical models are popularly used in computer vision, computational biology, digital communication, statistical physics, etc. Figure 1.1b illustrates an exemplified undirected graphical model in digital communication context. On the one hand, a receiver wants to know what is the true symbol by joint considering its communication environment and its received symbol. On the other hand, the symbol received by the receiver is jointly formulated by the true symbol and communication environment. The influence among them is apparently not directional since the impact along an edge can be bidirectional in this example.

Probabilistic graphical model offers a ‘scientific language’ to do reasoning with

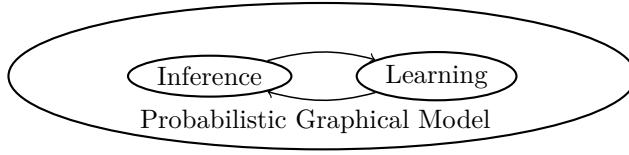


Figure 1.2: Two key aspects in practical graphical models.

uncertainty within framework of probabilistic theory. It is usually a nature representation for a complex system or problem and offers straightforward abstraction. The compact representation of probabilistic graphical model bridges the joint distribution of a complex system, and its graphical abstraction that captures the statistic dependency reflecting our understanding of the system. The advantage of its representation power is one of the reasons that leads to its popularity in difference disciplines.

Probabilistic graphical model coupled with its underlining distribution is a powerful tool for effective inference, apart from its advantage of representation power. It allows to answer queries with the help of the underlining distribution when practical inference algorithms are provided, which meets our need of reasoning with uncertainty. In addition to inference, probabilistic graphical model also supports learning from data. With certain amount of data available, a probabilistic graphical model can be learned to explain the observed data better in addition to align with our own understanding of a domain. The learned graphical model can serve to do inference with higher confidence in return. A diagram is illustrated in Figure 1.2. As would become clear in Part II, the inference may be needed to carry out model learning as well, apart from the above mutual-benefiting interaction.

1.2 Scope and Thesis Outline

We gives the intuition and motivation of probabilistic graphical model in last section, and the interaction between inference and learning in this framework. In this section, we would navigate further among the topics within this framework and states the ones that we would cover in the thesis.

Inference in probabilistic graphical model is about answering queries with help of its coupled distribution. These queries can be generally grouped into the following cases:

- Computing the likelihood of observed data or unobserved random variable.
- Computing the marginals distribution over a particular subset of random variables.
- Computing the conditional distribution a subset of variables given the configuration of another subset of variables.
- Computing the most likely configuration of (a subset of) variables.

The work of this thesis would be mainly related with the first three cases in inference part.

Due to either the requirement of efficiency in solving a problem or the structure of the problem's graphical model representation, it is not always that case that the above inference problem can be solved exactly. Thus inference methods can be divided into

- exact inference,
- and approximate inference.

For a limited class of graphs, exact inference such as variable elimination and sum-product algorithm can be used. Some graphs also allow efficient inference after mild graph modification, e.g. junction tree method. However, the above listed inference problems can only be approximately solved in general graphs. The approximate inference family can be further divided into

- Stochastic Approximation (Particle methods),
- Deterministic Approximation (Variational methods).

Stochastic approximation mainly relies on sample instances to answer queries, where a major challenge lies in how to obtain samples efficiently from a target distribution. Gibbs sampling, importance sampling and Markov Chain Monte Carlo are within this family. On the other hand, deterministic approximations refer to the variational methods, such as mean field approximation, loopy belief propagation, expectation propagation etc. *From the perspectives of methodology, we related work in this thesis locates in the family of variational methods under approximate inference category.*

As for learning in probabilistic graphical models, there are two types of learning problems

- Structure learning,
- Parameter learning.

The first case refer to determine the structure of a graphical model from observations, which is usually reduced to the problem of whether there should be an edge between a pair of nodes in the graphical model. The parameter learning is about to determine the parameter of a probabilistic graphical model (or its coupled distribution), with its graphical structure known. Structure learning is out of the scope of this thesis. The term *learning* in thesis means the estimation of the parameters of a distribution. This problem is mainly discussed in Part II, where we would touch the topics about learning in both undirected and directed graphical models.

As for the learning techniques, the learning principles can be categorized into

- Maximal likelihood estimation
- Maximal conditional likelihood

- Bayesian estimation
- Maximal ‘Margin’
- Maximum entropy

in general. We would touch techniques of the first four cases in Part II.

Publications

- Dong Liu and Lars K Rasmussen. Region-based energy neural network for approximate inference. Under review, 2020.
Code: <https://github.com/FirstHandScientist/renn> ([Private](#), [publish later](#))
- Dong Liu, Minh Thành Vu, Li Zuxing, and Lars K Rasmussen. α belief propagation for approximate bayesian inference. Under review, 2020.
Code: <https://github.com/FirstHandScientist/alpha-bp>
- Dong Liu, Antoine Honoré, Saikat Chatterjee, and Lars K Rasmussen. Powering hidden markov model by neural network based generative models. In the 24th European Conference on Artificial Intelligence (ECAI), 2020.
Code: <https://github.com/FirstHandScientist/genhmm>
- Dong Liu, Minh Thành Vu, Saikat Chatterjee, and Lars K Rasmussen. Neural network based explicit mixture models and expectation-maximization based learning. In International Joint Conference on Neural Networks, Glasgow, UK, July 2020.
Code: <https://github.com/FirstHandScientist/EM-GM>
- Antoine Honoré, Dong Liu, David Forsberg, Karen Coste, Eric Herlenius, Saikat Chatterjee, and Mikael Skoglund. Hidden markov models for sepsis detection in preterm infants. In 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2020.
- D. Liu, N. N. Moghadam, L. K. Rasmussen, J. Huang, and S. Chatterjee. α belief propagation as fully factorized approximation. In 2019 IEEE Global Conference on Signal and Information Processing (GlobalSIP), pages 1-5, 2019.
Code: <https://github.com/FirstHandScientist/amp>
- Dong Liu, Minh Thành Vu, Saikat Chatterjee, and Lars K Rasmussen. Entropy-regularized optimal transport generative models. In 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 3532-3536, 2019.
Code: <https://github.com/FirstHandScientist/eotgms>

- Dong Liu, Baptiste Cavarec, Lars K Rasmussen, and Jing Yue. On dominant interference in random networks and communication reliability. In ICC 2019 IEEE International Conference on Communications (ICC), pages 1-7, 2019.
- Dong Liu, Viktoria Fodor, and Lars K Rasmussen. Will scale-free popularity develop scale-free geo-social networks? IEEE Transactions on Network Science and Engineering, 2018.
- Dong Liu, Chao Wang, and Lars K Rasmussen. Discontinuous reception for multiple-beam communication. IEEE Access, pages 46931-46946, 2019.

Outline of Thesis

To be written after main content is finished.

Chapter 2

Background

In this chapter, we review some background knowledge that is going to be used in this thesis. We begin with the introduction to probabilistic graphical models. Then a divergence measure is introduced. Common inference tasks and methods are discussed before the learning problems in probabilistic graphical models are reviewed, which are interpreted as minimization of the divergence measure.

2.1 Graphical Models

Graphical models provide a formal graph representation of statistical dependency of complex problems or systems. The conditional independence of random variables can be conveniently encoded and analyzed by a graphical model. More importantly, query problems can be resolved by interactions of local regions of a graphical model in exact or approximate ways, which are usually unfeasible to solve directly.

More formally, a graphical model is a graphical representation of a collection of random variables (along their domains) where their statistical dependency is encoded into a set of non-negative functions and the graphical structure. Let $\mathbf{x} = (x_1, x_2, \dots, x_N)$ be a vector of random variables with N as a positive integer, where an element variable x_i can be either discrete or continuous random variable and takes values from its domain \mathcal{X}_i . Note that the domain of a random variable is not necessarily the same as that of another. With some abuse of notation, we might use \mathbf{x} to denote its assignment when there is no cause of ambiguity in context. The joint probability is denoted by $p(\mathbf{x}) = p(x_1, x_2, \dots, x_N)$. We denote $\mathcal{X} = \prod_{i=1}^N \mathcal{X}_i$ and then $\mathbf{x} \in \mathcal{X}$.

As motivated in Chapter 1.1, a graphical model can be directed or undirected. A directed graphical model is also known as a Bayesian network or generative model in literature [Bishop(2006), Chapter 8]. We might use the names alternatively. The non-negative functions in graphical models encode the local compatibility of states of random variables. In directed graphical models, i.e. Bayesian networks, the local functions are conditional probability functions. The joint prob-

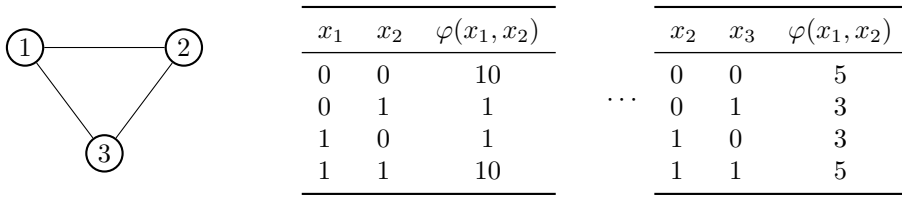


Figure 2.1: A Markov random field with three binary nodes. Potential factors are represented by tables.

ability distribution is represented as the product of these conditional probability functions,

$$p(\mathbf{x}) = \prod_{n=1}^N p(x_n | \mathcal{P}(x_n)), \quad (2.1)$$

where $\mathcal{P}(\cdot)$ denotes the set of parent nodes in the directed graph. In an directed graphical model, the local functions, i.e. the conditional probability distributions, e.g. $\{p(x_n | \mathcal{P}(x_n))\}$, are normalized and proper distributions. Additionally, sampling from a underlining distribution $p(\mathbf{x})$ of a directed graphical model is efficient. Due to acyclic property of directed graphical models, by the well know *ancestral sampling*, a sample (x_1, x_2, \dots, x_N) can be drawn sequentially via following the directed edges. In another word, x_n is always sampled after $\mathcal{P}(x_n)$. This process might be viewed as the 'generative' process of signal \mathbf{x} , i.e. how \mathbf{x} is generated from the graphical model.

A Bayesian network (generative model) is usually easier to be interpreted due to the fact that its local functions are conditional probabilities and it is natural to decompose the joint underlining distribution into conditional probability distributions. But Bayesian networks can only be applied to the limited cases where influence between variables is directional. In many practical cases, interaction between variables can not be naturally described by impact with directionality. A problem of this kind can be represented by an undirected graphical model, i.e. a Markov random field (MRF). Under certain condition, a Bayesian network can be perfectly represented by a Markov random field without loss of independence information by moralizing edges [Koller and Friedman(2009), Chapter 4.5]. Instead of conditional probability distributions, the local functions of MRF represents the compatibility of states of different variables, which are termed as *potential factors*. Different from conditional probabilities in a Bayesian network, a potential factor in a MRF is not necessary normalized (not necessary to be summed to one). We provide a toy example of MRF with three variable nodes as follows.

Example 2. As shown in Figure 2.1, the MRF encodes dependency of three random variables x_1 , x_2 , and x_3 , where node i is associated with variable x_i and each has a binary domain, i.e. $\mathcal{X}_i = \{0, 1\}$ for $i = 1, 2, 3$. Three potential factors of the MRF

together define the joint distribution

$$p(\mathbf{x}) = \frac{1}{Z} \varphi_{1,2}(x_1, x_2) \varphi_{2,3}(x_2, x_3) \varphi_{1,3}(x_1, x_3)$$

where $Z = \sum_{x_1, x_2, x_3} \varphi_{1,2}(x_1, x_2) \varphi_{2,3}(x_2, x_3) \varphi_{1,3}(x_1, x_3)$ normalizes the potential factors such that $p(\mathbf{x})$ is a proper distribution. The exemplified potential factors in Figure 2.1 demonstrate that it is more compatible or likely when x_1 , x_2 and x_3 are in the same state (either 0 or 1) than they are configured into different states.

From the above example to a formal statement, a MRF over random vector \mathbf{x} can be represented by a undirected graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$, with each node $i \in \mathcal{V}$ is associated with a random variable x_i and undirected edge set $\mathcal{E} \subset \mathcal{V} \times \mathcal{V}$. This MRF encodes a collection of distributions that factorize as

$$p(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \prod_{\alpha \in \mathcal{I}} \varphi_{\alpha}(\mathbf{x}_{\alpha}; \boldsymbol{\theta}), \quad (2.2)$$

where \mathcal{I} is the set of indexes of potential factors, and each factor φ_{α} for $\alpha \in \mathcal{I}$ is defined on subset of \mathbf{x} , i.e. $\varphi_{\alpha} : \mathcal{X}_{\alpha} \rightarrow \mathbb{R}^+ \cup \{0\}$, where $\mathcal{X}_{\alpha} = \prod_{i \in \alpha} \mathcal{X}_i$ is the domain of potential factor φ_{α} . The scope of factor α is $\mathbf{x}_{\alpha} = \{x_i | i \in \alpha\}$ where $i \in \alpha$ stands for that the variable x_i associated with node i is an argument of potential factor φ_{α} . In (2.2),

$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{x}} \prod_{\alpha \in \mathcal{I}} \varphi_{\alpha}(\mathbf{x}_{\alpha}; \boldsymbol{\theta}) \quad (2.3)$$

is the *partition function*. Apparently, the partition function normalizes the potential factors such that $p(\mathbf{x}; \boldsymbol{\theta})$ is a proper probability.

Remark 1. We can compare directed and undirected graphical models with regarding to the following aspects.

- *Representation:* The structure and the parameterization in directed graphical models provide a natural representation for many types of real-world domains. MRF representation is not usually as intuitive as that of directed graphical models. But the acyclic property of directed graphical models limits their representation power. On the other hand, MRFs can be either cyclic or acyclic, which offers the flexibility of graph structure and can simplify the graphical representation. Due to the weaker requirement of potential factors and the weaker requirement of graphical structure in MRFs than local functions and acyclic constraint in directed graphical models, respectively, the representation of MRFs are richer.
- *Local nonnegative functions:* The local functions are conditional probability functions in directed graphical models, but potential factors (nonnegative) in undirected cases.

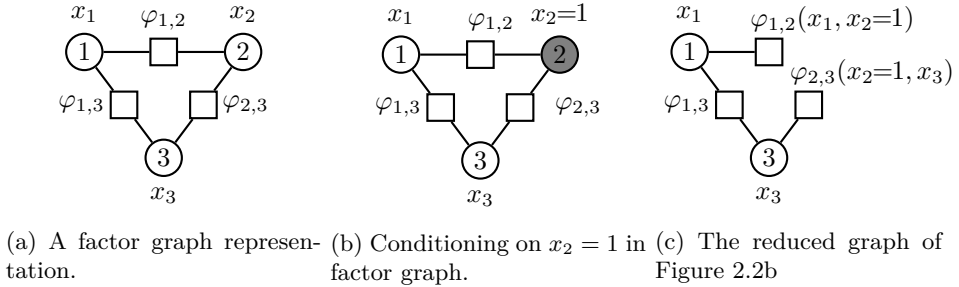


Figure 2.2: A Markov random field is represented by a factor graph, i.e. 2.2a, conditioning of the MRF 2.2b, the reduced MRF 2.2c.

- *Sampling:* Sampling is more straightforward within generative models (directed graphs) than that in MRFs.
- *Normalization:* Since each local function is a conditional probability function in directed graphical models, partition function for normalization is not needed. A MRF in general comes with partition functions, since potential factors are not necessarily normalized.

Alternative Representation of MRF

The representation of a MRF by $\mathcal{G}(\mathcal{V}, \mathcal{E})$ as explained above is compact, but the potential factors are not present in the graphical representation. An alternative representation to MRF is *factor graph* [Kschischang *et al.*(2001)Kschischang, Frey, and Loeliger], which is a bipartite graph topology. In a factor graph, a potential factor is explicitly represented as a factor node, as counterpart of variable node associated with a random variable.

Definition 1. A factor graph \mathcal{G}_F , is a bipartite graph that represents the factorization structure of (2.2). A factor graph has two types of nodes: i) a variable node for each variable x_i ; ii) a factor node for each potential function φ_α . An edge between a variable node i and factor node α if and only if x_i is argument of φ_α . We would denote a factor graph by $\mathcal{G}_F(\mathcal{V} \cup \mathcal{F}, \mathcal{E}_F)$ with \mathcal{V} as the set of variable nodes, \mathcal{F} as the set of factor nodes, and \mathcal{E}_F the set of undirected edges.

Example 3. Let us represent the Example 2.2a by a factor graph, which is shown in Figure 2.2a. Different from the representation by $\mathcal{G}(\mathcal{V}, \mathcal{E})$ in Figure 2.1, factor nodes are explicitly represented by square nodes.

Conditioning on Observations in MRFs

It is not rare that a graphical model may contain observed variable. The node set of a MRF can be separated into a subset \mathcal{V}_O of nodes, that are associated with

observed variable \mathbf{x}_O , and a subset \mathcal{V}_U of nodes associated with unobserved variable \mathbf{x}_U . For an evidence is observed,

$$p(\mathbf{x}_U|\mathbf{x}_O;\boldsymbol{\theta}) = \frac{p(\mathbf{x}_U, \mathbf{x}_O;\boldsymbol{\theta})}{p(\mathbf{x}_O;\boldsymbol{\theta})} = \frac{Z(\mathbf{x}_O, \boldsymbol{\theta})}{Z(\boldsymbol{\theta})}, \quad (2.4)$$

where

$$\begin{aligned} \tilde{p}(\mathbf{x};\boldsymbol{\theta}) &= \prod_{\alpha \in \mathcal{I}} \varphi_{\alpha}(\mathbf{x}_{\alpha};\boldsymbol{\theta}), \\ Z(\mathbf{x}_O, \boldsymbol{\theta}) &= \sum_{\mathbf{x}_U} \tilde{p}(\mathbf{x};\boldsymbol{\theta}), \\ Z(\boldsymbol{\theta}) &= \sum_{\mathbf{x}_O} \sum_{\mathbf{x}_U} \tilde{p}(\mathbf{x};\boldsymbol{\theta}). \end{aligned} \quad (2.5)$$

This means that a condition probability can be computed by partition function and sub-partition functions. Alternatively, when an evidence \mathbf{e}_O (an sample instance of \mathbf{x}_O) is observed, the conditional probability can be written as

$$p(\mathbf{x}_U|\mathbf{x}_O = \mathbf{e}_O;\boldsymbol{\theta}) = \frac{\tilde{p}(\mathbf{x}_U, \mathbf{x}_O = \mathbf{e}_O;\boldsymbol{\theta})}{\sum_{\mathbf{x}_U} \tilde{p}(\mathbf{x}_U, \mathbf{x}_O = \mathbf{e}_O;\boldsymbol{\theta})} \propto \tilde{p}(\mathbf{x}_U, \mathbf{x}_O = \mathbf{e}_O;\boldsymbol{\theta}) \quad (2.6)$$

where \propto stands for propositional to. 2.6 shows an interesting phenomenon for MRF including evidence. It can be understood as clamping nodes in \mathcal{V}_O of the MRF to configuration \mathbf{e}_O , i.e. the domain of \mathbf{x}_O becomes a set containing only one instance \mathbf{e}_O . For instance, an example of conditioning on a variable node for Example 2 is shown in Figure 2.2b.

In addition to the above intuitions, conditioning can also be understood as a process of reducing the graph of a MRF. When a MRF is conditioned on \mathbf{x}_O , the variables nodes of set \mathcal{V}_O are removed from \mathcal{G} , along with their edges. The potential factors with regarding to \mathcal{V}_U are modified accordingly [Koller and Friedman(2009), Chapter 4.2.3]. For instance, the graph including evidence node 2 in Figure 2.2b can be further reduced into a Figure 2.2c. Then any inference applicable to a MRF applies to the MRF with nodes clamped as well. A MRF with several nodes clamped to some evidence can be seen either as a manipulation of its domain or the graph itself.

It can be seen that MRF framework is capable to handle conditioning as well. Therefore, in the following part of the thesis, it might or might not have been based on conditioning observed variables when a MRF is mentioned.

2.2 Divergence

Before we get into more discussion about inference and learning topics, we firstly introduce the concept of *divergence* measures since principles of both learning and

inference are closely related with divergence measure. A divergence measure plays a fundamental role when we try to use a probability distribution (over discrete or continuous variable) q to approximate another probability distribution p . A divergence measure is used to formally quantify how much information is lost when p is represented by q . Denote \mathbb{P} as the space of measures p and q , i.e. $p, q \in \mathbb{P}$.

Definition 2. *Given the space \mathbb{P} of probability distribution for a random variable \mathbf{x} , a divergence on this space is defined as a function $D(p||q) : \mathbb{P} \times \mathbb{P} \rightarrow \mathbb{R}^+ \cup \{0\}$ such that $D(p||q) \geq 0$ for all $p, q \in \mathbb{P}$ and $D(p||q) = 0$ if and only if $p = q$.*

Here we introduce the classic *Kullback-Leibler divergence* [Kullback(1959), Kullback and Leibler(1951)], KL divergence for short, which is one of the most widely used divergence measures in machine learning, statistics and information theory.

Definition 3. *The Kullback-Leibler (KL) divergence on \mathbb{P} is defined as a function $KL(\cdot||\cdot) : \mathbb{P} \times \mathbb{P} \rightarrow \mathbb{R}^+ \cup 0$ with the following form*

$$KL(p||q) = \sum_{\mathbf{x}} p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})}, \quad (2.7)$$

where \log is the natural logarithm. Note the sum in 2.7 should be replaced by integral when p and q are probability density functions.

KL divergence is not symmetric. In another word, there is no equivalence between $KL(p||q)$ and $KL(q||p)$ in general.

2.3 Inference Tasks

Given a probability distribution $p(\mathbf{x})$ as the underline distribution of a graphical model, inference in general can be divided into four kinds of tasks, as brought up in Chapter 1.2. Our work in this thesis would be closely involved with the problems

- computing the likelihood of observed data or unobserved random variable;
- computing the marginal distribution over a particular subset of nodes, i.e. $p(\mathbf{x}_A)$ for $A \in \mathcal{V}$. Note that single-node marginal distribution $p(x_i)$ also belongs to this case;
- computing the conditional distribution a subset of nodes given the configuration of another subset of nodes, i.e. $p(\mathbf{x}_A|\mathbf{x}_B)$ for $A, B \in \mathcal{V}$ and $A \cap B = \emptyset$;

in MRFs. The above tasks are also closely related with the inference of partition function

- Computation of $Z(\boldsymbol{\theta}) = \sum_{\mathbf{x}} \prod_{\alpha \in \mathcal{I}} \varphi_{\alpha}(\mathbf{x}_{\alpha}; \boldsymbol{\theta})$, or sub-partition functions.

In the following section, we would introduce the variational methods for the above tasks in high-level.

2.4 Variational inference

In solving inference tasks, one important technique is based on a variational approach. With $p(\mathbf{x}; \boldsymbol{\theta})$ as the underlining probability distribution of a graphical model, directly inference with $p(\mathbf{x}; \boldsymbol{\theta})$ is often unfeasible due to the system represented by the graphical model is too large or complex. It can also be the case that even we know the form of $p(\mathbf{x}; \boldsymbol{\theta})$, the computation in inference tasks can be prohibitive. In variational approaches, a 'trial' probability distribution $b(\mathbf{x})$ is introduced to approximate $p(\mathbf{x}; \boldsymbol{\theta})$. The trial distribution should be intuitively simpler than $p(\mathbf{x}; \boldsymbol{\theta})$. *Variational free energy* [Oppen and Saad(2001)] is a quantity used to find such a approximation. The variational free energy is defined by

$$\begin{aligned} F_V(b) &= \text{KL}(b(\mathbf{x})||p(\mathbf{x}; \boldsymbol{\theta})) - \log Z(\boldsymbol{\theta}) \\ &= \sum_{\mathbf{x}} b(\mathbf{x}) \log \frac{b(\mathbf{x})}{p(\mathbf{x}; \boldsymbol{\theta})} - \log Z(\boldsymbol{\theta}) \\ &= \sum_{\mathbf{x}} b(\mathbf{x}) \log \frac{b(\mathbf{x})}{\tilde{p}(\mathbf{x}; \boldsymbol{\theta})}, \end{aligned} \quad (2.8)$$

where $\tilde{p}(\mathbf{x}; \boldsymbol{\theta}) = \prod_{\alpha \in \mathcal{I}} \psi_{\alpha}(\mathbf{x}_{\alpha}; \boldsymbol{\theta}_{\alpha})$. Since $\text{KL}(b(\mathbf{x})||p(\mathbf{x}; \boldsymbol{\theta}))$ is always non-negative and is zero if and only if $b(\mathbf{x}) = p(\mathbf{x}; \boldsymbol{\theta})$, we have $F_V(b) \geq -\log Z(\boldsymbol{\theta})$, with equality when $b(\mathbf{x}) = p(\mathbf{x}; \boldsymbol{\theta})$.

Remark 2. Note from 2.8, the minimization w.r.t. b of variational free energy is equivalent to the divergence minimization, i.e. $\text{KL}(b(\mathbf{x})||p(\mathbf{x}; \boldsymbol{\theta}))$, since $\log Z(\boldsymbol{\theta})$ does not depend on b . By observing $F_V(b) = \sum_{\mathbf{x}} b(\mathbf{x}) \log \frac{b(\mathbf{x})}{\tilde{p}(\mathbf{x}; \boldsymbol{\theta})}$, the free energy minimization guides the choice of b such that b is close to an unnormalized measure $\tilde{p}(\mathbf{x}; \boldsymbol{\theta})$ in its space.

Another benefit of 2.8 is that we are able to approximate $p(\mathbf{x}; \boldsymbol{\theta})$ without inference of the true marginal distributions $\{p(\mathbf{x}_{\alpha}; \boldsymbol{\theta}), \alpha \in \mathcal{I}\}$. Since $\log \tilde{p}(\mathbf{x}; \boldsymbol{\theta})$ can be formulated as sum of log-potential-factors that are all local functions, the computation of $F_V(b)$ can be done by inference of marginals $\{b_{\alpha}(\mathbf{x}_{\alpha}), \alpha \in \mathcal{I}\}$ of the approximate distribution, which are tractable.

Remark 3. Discussion: Since we are essentially approximating distribution p by a distribution b , can we minimize $\text{KL}(p(\mathbf{x}, \boldsymbol{\theta})||b(\mathbf{x}))$ instead?

It might be feasible by instinct. But a further check would reveal its infeasibility. The $\text{KL}(p(\mathbf{x}; \boldsymbol{\theta})||b(\mathbf{x}))$ would inevitable requires the marginals of p and therefore requires the exact inference in p , which are what we are trying to avoid. But it does not mean this divergence is useless. As shall be seen in section 2.5, this type of divergence measure is what we need in model learning.

Variational Free Energy and Mean Field

In mean field approach, a fully-factorized approximation is used,

$$b_{MF}(\mathbf{x}) = \prod_{i=1}^N b_i(x_i). \quad (2.9)$$

Substituting (2.9) into the variational free energy gives

$$F_{MF}(b) = - \sum_{\alpha \in \mathcal{I}} \sum_{\mathbf{x}_\alpha} \log \varphi_\alpha(\mathbf{x}_\alpha; \boldsymbol{\theta}) \prod_{i \in \alpha} b_i(x_i) + \sum_{i \in \mathcal{V}} \sum_{x_i} b_i(x_i) \log b_i(x_i), \quad (2.10)$$

where $i \in \alpha$ stands for the node i 's associated x_i is argument of φ_α . Solving the minimization of F_{MF} w.r.t. $b_{MF}(\mathbf{x})$ gives the update rule of mean field

$$\log b_i(x_i) \propto \sum_{\alpha \in \text{ne}_i} \sum_{\mathbf{x}_\alpha \setminus x_i} \log \varphi_\alpha(\mathbf{x}_\alpha; \boldsymbol{\theta}_\alpha) \prod_{j \in \alpha \setminus i} b_j(x_j), \quad (2.11)$$

where $\text{ne}_i = \{\alpha | i \in \alpha, \alpha \in \mathcal{F}\}$, i.e. the potential factors that has x_i as argument.

Bethe Free Energy and (Loopy) Belief Propagation

Different from the mean field approximation, Bethe approximation also includes the non-single-node beliefs $\{b_\alpha(\mathbf{x}_\alpha)\}$ apart from the single-node beliefs $\{b_i(x_i)\}$ [Yedidia *et al.* (2003) Yedidia, Freeman, and Weiss]. In this case, the Bethe free energy is given by

$$F_{Bethe}(b) = \sum_{\alpha} \sum_{\mathbf{x}_\alpha} b_\alpha(\mathbf{x}_\alpha) \log \frac{b_\alpha(\mathbf{x}_\alpha)}{\varphi_\alpha(\mathbf{x}_\alpha)} - \sum_{i=1}^N (|\text{ne}_i| - 1) \sum_{x_i} b_i(x_i) \log b_i(x_i), \quad (2.12)$$

where $|\cdot|$ stands for cardinality. Due to the non-single-node beliefs, there are consistency constraints $\sum_{\mathbf{x}_\alpha} b_\alpha(\mathbf{x}_\alpha) = \sum_{x_i} b_i(x_i) = 1, \forall i \in \alpha$ to obey. Then, solving the Bethe free energy minimization problem

$$\begin{aligned} & \min_{\{b_\alpha(\mathbf{x}_\alpha)\}, \{b_i(x_i)\}} F_{Bethe}(b) \\ & \text{s.t.} \quad \sum_{\mathbf{x}_\alpha \setminus x_i} b_\alpha(\mathbf{x}_\alpha) = b_i(x_i), \\ & \quad \sum_{\mathbf{x}_\alpha} b_\alpha(\mathbf{x}_\alpha) = \sum_{x_i} b_i(x_i) = 1, \\ & \quad 0 \leq b_i(x_i) \leq 1, \\ & \quad i \in \mathcal{V}, \alpha \in \mathcal{F}, \end{aligned} \quad (2.13)$$

where \mathcal{V} and \mathcal{F} are the set of variable nodes and the set of factor nodes in factor graph as defined in Definition 1, gives the (loopy) BP message-passing rule

$$m_{a \rightarrow i}(x_i) \propto \sum_{\mathbf{x}_\alpha \setminus x_i} \varphi_\alpha(\mathbf{x}_\alpha) \prod_{j \in \alpha \setminus i} \prod_{\alpha' \in \text{ne}_j \setminus \alpha} m_{\alpha' \rightarrow j}(x_j). \quad (2.14)$$

2.5 Learning principles

We have touched the learning topic in chapter 1, which is to find the 'best' probability distribution $p(\mathbf{x}; \boldsymbol{\theta})$ in its space \mathbb{P} . To make the discussion more concrete, we assume the domain is governed by a underlying distribution p^* that is induced by a (directed or undirected) graphical model, $\mathcal{M}^* = \{\mathcal{K}^*, \boldsymbol{\theta}^*\}$ with \mathcal{M}^* representing its structure and $\boldsymbol{\theta}^*$ representing its parameter. Here we discuss about *model learning* (parameter learning only). For notation simplicity, we use $p^*(\mathbf{x})$ to denote this distribution. We are given a dataset $\mathcal{D} = \{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^M\}$. Following the standard assumption, these sample instances are *independent and identically distributed (i.i.d.)*. The task is then to use the information from the dataset to learn a distribution p within its space \mathbb{P} , since the governing distribution $p^*(\mathbf{x})$ is not known.

The problem of learning a distribution in \mathbb{P} to approximate p^* can be formulated as density estimation. With the concept of KL divergence in section 2.2, learning of p can be formulated as minimizing the KL divergence

$$\begin{aligned} & \text{KL}(p^*(\mathbf{x}) \| p(\mathbf{x}; \boldsymbol{\theta})) \\ &= \mathbb{E}_{\mathbf{x} \sim p^*} \left[\log \frac{p^*(\mathbf{x})}{p(\mathbf{x}; \boldsymbol{\theta})} \right] \\ &= -H(p^*) - \mathbb{E}_{\mathbf{x} \sim p^*} [\log p(\mathbf{x}; \boldsymbol{\theta})], \end{aligned} \quad (2.15)$$

where $H(p^*)$ is the entropy of p^* . Due to the property of divergence, the KL divergence in 2.15 is zero if and only if $p(\mathbf{x}; \boldsymbol{\theta}) = p^*(\mathbf{x})$. The last line of 2.15 shows that the negative entropy term does not depends on $p(\mathbf{x}; \boldsymbol{\theta})$. Thus we can just focus on the expectation term $\mathbb{E}_{\mathbf{x} \sim p^*} [\log p(\mathbf{x}; \boldsymbol{\theta})]$, which is *expected log-likelihood*. Therefore, we can just use the expected log-likelihood to do model learning instead of minimizing the KL divergence.

Note although we can use the expected log-likelihood for model learning task and even model comparison (comparing a trained model with another one), we loss the information of how close a trained model is to p^* . This is due to the omitting of $H(p^*)$, which is not available.

Since it is not possible to know p^* (otherwise we do not need to learn it), the expected log-likelihood is approximated by sample instances of p^* ,

$$\mathcal{L}(\mathcal{D}; \boldsymbol{\theta}) = \frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \log p(\mathbf{x}; \boldsymbol{\theta}), \quad (2.16)$$

and

$$\mathbb{E}_{\mathbf{x} \sim p^*} (\log p(\mathbf{x}; \boldsymbol{\theta})) \approx \mathcal{L}(\mathcal{D}; \boldsymbol{\theta}). \quad (2.17)$$

Log-likelihood $\mathcal{L}(\mathcal{D}; \boldsymbol{\theta})$ is one of the most widely used loss for model learning. However, $\mathcal{L}(\mathcal{D}; \boldsymbol{\theta})$ is not always a feasible loss to compute due to:

- exact computation of $p(\mathbf{x})$ is not possible;

- there are some elements of \mathbf{x} which are not observable (latent variables).

For the first case, the typical treatment is to approximate the exact log-likelihood. This is done by approximation with employing inference methods or making simplified assumptions on dependency structure of the graphical model of $p(\mathbf{x})$. Then, optimization is carried out with regarding to the approximated log-likelihood. These methods include surrogate likelihood [Wainwright(2006), Lu *et al.*(2018) Lu, Liu, and Huang], pseudo-likelihood [Qu *et al.*(2019) Qu, Bengio, and Tang, Lazaro-Gredilla *et al.*(2019) Lazaro-Gredilla, Lehrach, and George], piecewise likelihood [Sutton and McCallum(2012), Lin *et al.*(2016) Lin, Shen, van den Hengel, and Reid], saddle-point approximation [Srikumar *et al.*(2012) Srikumar, Kundu, and Roth, Wiseman and Kim(2019)].

Apart from the above case where all variables are observable, the partial observed models, the latent variable case, are equally import in inference and learning with uncertainty. This class of models includes (but not limited to) classic Gaussian mixture models (GMMs) and hidden Markov models (HMMs). There are latent variables because:

- Use of abstract variable to model the generative process (usually a directed graph) of observation data, such as HMMs.
- A practical true attribute of an object may be difficult or impossible to measure exactly. For instance, the disease infection can only be diagnosed via the relevant symptoms, e.g. Example 1; In the position tracking of a car with noisy sensors, the true position of the car might only be inferred via noisy data of sensors.
- No measurement on an attribute of an object of interest is made. For instance, the velocity sensor in the car tracking example might not be stable and might read not quantity now and then.

In general, latent variables are commonly used to deal with partial observation problems, data clustering, data manipulation, etc. Let us denote the observable variable and latent variable by \mathbf{x}_O and \mathbf{x}_U , respectively. We can see that the log-likelihood $p(\mathbf{x}_O, \mathbf{x}_U; \boldsymbol{\theta})$ is not available any more as it is in the fully-observed case. To deal with the latent variables, we can try to optimize the partial log-likelihood

$$\begin{aligned}
 l(\mathbf{x}_O; \boldsymbol{\theta}) &= \log p(\mathbf{x}_O; \boldsymbol{\theta}) \\
 &= \mathbb{E}_{q(\mathbf{x}_U|\mathbf{x}_O)} \left[\log \frac{q(\mathbf{x}_U|\mathbf{x}_O)}{p(\mathbf{x}_U|\mathbf{x}_O; \boldsymbol{\theta})} \cdot \frac{p(\mathbf{x}_U, \mathbf{x}_O; \boldsymbol{\theta})}{q(\mathbf{x}_U|\mathbf{x}_O)} \right] \\
 &= \text{KL}(q(\mathbf{x}_U|\mathbf{x}_O) \| p(\mathbf{x}_U|\mathbf{x}_O; \boldsymbol{\theta})) + F(q, \boldsymbol{\theta})
 \end{aligned} \tag{2.18}$$

with

$$F(q, \boldsymbol{\theta}) = \mathbb{E}_{q(\mathbf{x}_U|\mathbf{x}_O)} [\log p(\mathbf{x}_U, \mathbf{x}_O; \boldsymbol{\theta})] + H(q(\mathbf{x}_U|\mathbf{x}_O)) \tag{2.19}$$

where $H(q(\mathbf{x}_U|\mathbf{x}_O))$ is the entropy of $q(\mathbf{x}_U|\mathbf{x}_O)$, and q can be any distribution over \mathbf{x}_U . Due to the non-negative property of KL divergence, we have

$$l(\mathbf{x}_O; \boldsymbol{\theta}) \geq F(q, \boldsymbol{\theta}), \quad (2.20)$$

with equality when $q(\mathbf{x}_U|\mathbf{x}_O) = p(\mathbf{x}_U|\mathbf{x}_O; \boldsymbol{\theta})$. $F(q, \boldsymbol{\theta})$ is also called *variational lower bound*.

One of the most wide used methods in learning with latent variable is *expectation maximization (EM)* [Dempster *et al.*(1977) Dempster, Laird, and Rubin]. In EM method, the posterior of \mathbf{x}_U is computed exactly from p , $q(\mathbf{x}_U|\mathbf{x}_O) = \text{argmax}_q F(q, \boldsymbol{\theta}) = p(\mathbf{x}_U|\mathbf{x}_O; \boldsymbol{\theta})$, which is the optimal solution to q . Then the parameter $\boldsymbol{\theta}$ of p is optimized. The two steps are optimized iteratively.

In cases the posterior of \mathbf{x}_U is not feasible to compute, the variational EM [Wainwright and Jordan(2008), section 6.2.2] or Monte Carlo EM (need sampling technique) [Neath(2012)]. There are also neural network based methods with Monte Carlo estimator to cope with the latent variable problems, see [Kingma and Welling(2014)] [Kuleshov and Ermon(2017), Lazaro-Gredilla *et al.*(2019) Lazaro-Gredilla, Lehrach, and George, Goodfellow *et al.*(2014) Goodfellow, Pouget-Abadie, Mirza, Xu, Warde-Farley, Ozair, Courville, and Bengio].

In above discussion, we have assumed that the model's distribution q is explicitly defined, i.e. the distribution (along its density function or mass function) is available. In an alternative track, the deep generative models grow popularity in literature and can be applied to different areas such as high-dimensional data representation, reinforcement learning and semi-supervised learning for its efficient sampling of multi-mode distributions [Goodfellow(2017)]. A generator in a deep generative model induces a distribution that can be either explicit (cite RealNVP, differential flow, autoencoders) or implicit (GAN) distribution. The former problem dates back to [Deco and Brauer(1995)] and receives more attention in recent years with latest work such as variational autoencoder [Kingma and Welling(2014)], Real-NVP [Dinh *et al.*(2016) Dinh, Sohl-Dickstein, and Bengio], normalizing flow model [Kingma and Dhariwal(2018)] and neural ordinary difference equations [Chen *et al.*(2018) Chen, Rubanova, Bettencourt, and Duvenaud]. The training of these models are still based on the (variational) maximum likelihood principle.

The later case brings a implicit distribution, where the maximum likelihood principle is not applicable any more. In this case, a state-of-art method is the generative adversarial model that employs a discriminator to play the role of divergence measure [Goodfellow *et al.*(2014) Goodfellow, Pouget-Abadie, Mirza, Xu, Warde-Farley, Ozair, Courville, and Bengio, Arjovsky and Bottou(2017)] [Salimans *et al.*(2016) Salimans, Goodfellow, Zaremba, *et al.*], which is essentially explained by a process of the minimization of the Jensen-Shannon divergences. Additionally, other sample based distances are employed as alternative methods for implicit model learning. Optimal transport is receiving more attention in recent years [Santambrogio(2015)], [Cuturi(2013)] in this track, which also applied for training of Boltzmann machine [Montavon *et al.*(2016) Montavon, Müller, and Cuturi] auto-encoders [Tolstikhin *et al.*(2017) Tolstikhin, Bousquet, Gelly, and Schoelkopf] and

generative adversarial networks [Arjovsky *et al.*(2017)Arjovsky, Chintala, and Bottou].

Part I

Inference

Chapter 3

An alternative view of belief propagation

Content:

1. α Belief Propagation as Fully Factorized Approximation, GlobalSIP 2019.
2. α Belief Propagation for Approximate Bayesian Inference, under review.

3.1 α belief propagation

3.2 Convergence study

3.3 Experimental results

3.4 Summary

Chapter 4

Region-based Energy Neural Network Model

work in Region-based Energy Neural Network for Approximate Inference, under, review

4.1 Region-based graph and energy

4.2 RENN model for Approximate Inference

4.3 RENN model for markov random field training

4.4 Experimental results

4.5 Summary

Part II

Learning

Chapter 5

Learning with inference

5.1 learning Undirected graphical models/ MRF

move the MRF learning by using RENN here

I should read lecture note 7 of 10-708 again when writing this section.

5.2 Amortized/Neural Variational Learning and Inference of partial observed MRF

1. TRW as upper bound to partition function
 2. Mean field or negative TRW as lower bound to partition function
- combining above together, we can obtain two different lower bound of likelihood.
Consider if worthy a paper.

- The log-likelihood of partial observed MRF is non-convex in general (log-sum-exp is convex, but the difference of two log-sum-exp functions might not be). This combination convert the original non-convex learning into convex optimization with regarding to MRF parameter? should be, but need a confirmation.
- 1. The speed of training can be improved by directly optimizing amortized beliefs.
- The bound becomes tighter by using clamping of variable, clamping can be done with or without selection of variables. No sampling is needed in training or inference.
- If need more contribution, use tree-reweighted hyper graph to obtain tighter bound.
- Not necessarily done here: the bound can also be further improved by important sampling.

Reference:

- 1. Wainwright, 2003, Tree-reweighted belief propagation algorithms and approximate ML estimation by pseudo-moment matching
- 2. Weller, 2015, Clamping Improves TRW and Mean Field Approximations
- 3. Mnih, 2014, Neural Variational Inference and Learning in Belief Networks, which describes a neural variational method for belief network. The major difference is the belief network as a DAG do not have the problem of partition function difficulty as MRF or partial observed MRF.

5.3 Notation

Random variable $\mathbf{v} \in \mathcal{X}_v$ that can be observed. Random variable $\mathbf{h} \in \mathcal{X}_h$ that is hidden variable and can not be observed.

An alternative plan:

- Training RENN with marginal-likelihood instead of joint likelihood, ref to Domke13
- If the above works, use Gaussian kernels to define potential, Marvin T. T. Teichmann Convolutional CRFs for Semantic Segmentation

5.4 Model and Problem Definition

We define the conditional probabilistic model as

$$p(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \tilde{p}(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta}), \quad (5.1)$$

with

$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{v}} \sum_{\mathbf{h}} \tilde{p}(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta}) \quad (5.2)$$

$$\tilde{p}(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta}) = \exp \{-E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta})\} \quad (5.3)$$

where $E(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta})$ is the average energy: $\mathcal{X}_v \times \mathcal{X}_h \rightarrow \mathbb{R}$.

We want to maximize the marginal likelihood:

$$\max_{\boldsymbol{\theta}} \log \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta}) = \max_{\boldsymbol{\theta}} \log Z(\mathbf{v}, \boldsymbol{\theta}) - \log Z(\boldsymbol{\theta}), \quad (5.4)$$

where $Z(\mathbf{v}, \boldsymbol{\theta}) = \sum_{\mathbf{h}} \tilde{p}(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta})$

5.5 A lower bound of the marginal likelihood

Denote $A(\boldsymbol{\theta}) = \log Z(\boldsymbol{\theta})$ and $A(\mathbf{v}, \boldsymbol{\theta}) = \log Z(\mathbf{v}, \boldsymbol{\theta})$

$$E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}) = -\langle \boldsymbol{\theta}, \boldsymbol{\varphi}(\mathbf{v}, \mathbf{h}) \rangle \quad (5.5)$$

and

$$\boldsymbol{\mu} = \mathbb{E}_{p(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta})}[\boldsymbol{\varphi}(\mathbf{v}, \mathbf{h})]. \quad (5.6)$$

In case of overcomplete representation of $\boldsymbol{\varphi}$, $\boldsymbol{\mu}$ is the set of marginal distributions.

With mean field approximation,

$$A_M(\mathbf{v}, \boldsymbol{\theta}) = \max_{\boldsymbol{\mu}_v \in \mathcal{M}_M} \langle \boldsymbol{\theta}, \boldsymbol{\mu}_v \rangle + H_M(\boldsymbol{\mu}_v), \quad (5.7)$$

where \mathcal{M}_M is the subspace of distributions where each variable is independent. And we have

$$A_M(\mathbf{v}, \boldsymbol{\theta}) \leq A(\mathbf{v}, \boldsymbol{\theta}). \quad (5.8)$$

With tree-reweighted approximation, TRW,

$$A_T(\boldsymbol{\theta}) = \max_{\boldsymbol{\mu} \in \mathcal{M}_T} \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle + H(\boldsymbol{\mu}), \quad (5.9)$$

where \mathcal{M} is the subspace of distributions where each variable is independent. And we have

$$A_T(\boldsymbol{\theta}) \geq A(\boldsymbol{\theta}). \quad (5.10)$$

We define the lower bound of marginal loglikelihood:

$$\mathcal{L}(\boldsymbol{\theta}) = A_M(\mathbf{v}, \boldsymbol{\theta}) - A_T(\boldsymbol{\theta}) \leq \log \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta}). \quad (5.11)$$

Connection to RBM:

$$p(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp\{\mathbf{v} \mathbf{W} \mathbf{h} + \mathbf{v} \mathbf{b} + \mathbf{v} \mathbf{a}\} \quad (5.12)$$

Note $p(\mathbf{h}|\mathbf{h})$ is exactly independent, and thus the $A_M(\boldsymbol{\theta}) = A(\boldsymbol{\theta})$ can be achieved, then how tight the lower bound $\mathcal{L}(\boldsymbol{\theta})$ would depend only on the TRW bound.

I should also consider how to use the trained model for prediction.

This is also closely connected to variational see Section 6.2, Wainwright, Graphical Models, Exponential Families, and Variational Inference.

5.6 Experiment

- Start with standard RBM section 4.2 in Amortized learning of MRFs
 - Try to break the conditional independence by connecting nodes of \mathbf{h}
 - Extend to conditional RBM training for denoising and data completion
- high-order HMMs

Chapter 6

Powering the expectation maximization method by neural networks

content: Neural Network based Explicit Mixture Models and Expectation-maximization based Learning, under review

section/chapter transition text: mixture model could be obtained from clamping and condition on a discrete variable, ref to Geier, Locally conditional belief propagation. Weller, clamping variables and approximate inference

Remark 4. *Theory interpretation of EM and variational EM, see Section 6.2, Wainwright, Graphical Models, Exponential Families, and Variational Inference.*

Remark 5. *RELATIONSHIP TO K-MEANS CLUSTERING Big picture: The EM algorithm for mixtures of Gaussians is like a soft version of the K-means algorithm.*

Remark 6. *EM lower bound + entropy of posterior of latent variable if a free energy. ref to 10-708 lecture6 note. EM using posterior of latent variable is equivalent to fully observable MLE where statistics are replaced by their expectations w.r.t the posterior.*

Can be viewed as two-node graphical model learning. 10-708lecture5-note

- 6.1** Normalizing flow
- 6.2** expectation maximization of neural network based mixture models
- 6.3** An alternative construction method
- 6.4** Experiments
- 6.5** Summary

Chapter 7

Powering Hidden Markov Model by Neural Network based Generative Models

For permeable part and notation of this chapter, refer to [Koller and Friedman(2009), Chapter 6.2]. Give a figure/illustration: Dynamic Bayesian Network > 2-TBN > HMM

A bit history of HMM, see [Koller and Friedman(2009), Chapter 6.8]
content:

1. Powering Hidden Markov Model by Neural Network based Generative Models, ECAI 2020

2. Antoine Honore, Dong Liu, Hidden Markov Models for sepsis detection in preterm infants, ICASSP, 2020

HMM is an instance of 2-time-slice Bayesian network(2-TBN) (section 6.2.2 Koller). Also, it can be argued from CRF.

7.1 Hidden Markov Model

7.2 GenHMM

7.3 Application to phone recognition

7.4 Application to sepsis detection in preterm infants

7.5 Summary

Chapter 8

An implicit probabilistic generative model

content: Entropy-regularized Optimal Transport Generative Models, ICASSP 2019

- 8.1 Modeling data without explicit probabilistic distribution
- 8.2 Employing EOT for modeling
- 8.3 Experimental results
- 8.4 Summary

Part III

Epilogue

Chapter 9

Conclusion and Discussions

Bibliography

- [Arjovsky and Bottou(2017)] M. Arjovsky and L. Bottou. 2017. Towards Principled Methods for Training Generative Adversarial Networks. *ArXiv e-prints*.
- [Arjovsky *et al.*(2017)Arjovsky, Chintala, and Bottou] Martin Arjovsky, Soumith Chintala, and Léon Bottou. 2017. Wasserstein generative adversarial networks. In *International Conference on Machine Learning*, pages 214–223.
- [Bishop(2006)] Christopher M. Bishop. 2006. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag, Berlin, Heidelberg. ISBN 0387310738.
- [Chen *et al.*(2018)Chen, Rubanova, Bettencourt, and Duvenaud] Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. 2018. Neural ordinary differential equations. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 6571–6583. Curran Associates, Inc. URL <http://papers.nips.cc/paper/7892-neural-ordinary-differential-equations.pdf>.
- [Cuturi(2013)] Marco Cuturi. 2013. Sinkhorn distances: Lightspeed computation of optimal transport. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 2292–2300. Curran Associates, Inc.
- [Deco and Brauer(1995)] Gustavo Deco and Wilfried Brauer. 1995. Higher order statistical decorrelation without information loss. In G. Tesauro, D. S. Touretzky, and T. K. Leen, editors, *Advances in Neural Information Processing Systems 7*, pages 247–254. MIT Press. URL <http://papers.nips.cc/paper/901-higher-order-statistical-decorrelation-without-information-loss.pdf>.
- [Dempster *et al.*(1977)Dempster, Laird, and Rubin] A. P. Dempster, N. M. Laird, and D. B. Rubin. 1977. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society: Series B*, 39:1–38. URL <http://web.mit.edu/6.435/www/Dempster77.pdf>.

- [Dinh *et al.*(2016)] Dinh, Sohl-Dickstein, and Bengio] L. Dinh, J. Sohl-Dickstein, and S. Bengio. 2016. Density estimation using Real NVP. *ArXiv e-prints*.
- [Goodfellow(2017)] I. Goodfellow. 2017. NIPS 2016 Tutorial: Generative Adversarial Networks. *ArXiv e-prints*.
- [Goodfellow *et al.*(2014)] Goodfellow, Pouget-Abadie, Mirza, Xu, Warde-Farley, Ozair, Courville, Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 2672–2680. Curran Associates, Inc. URL <http://papers.nips.cc/paper/5423-generative-adversarial-nets.pdf>.
- [Kingma and Dhariwal(2018)] D. P. Kingma and P. Dhariwal. 2018. Glow: Generative Flow with Invertible 1x1 Convolutions. *ArXiv e-prints*.
- [Kingma and Welling(2014)] Diederik P. Kingma and Max Welling. 2014. Auto-encoding variational bayes. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*.
- [Koller and Friedman(2009)] Daphne Koller and Nir Friedman. 2009. *Probabilistic Graphical Models: Principles and Techniques - Adaptive Computation and Machine Learning*. The MIT Press. ISBN 0262013193.
- [Kschischang *et al.*(2001)] Kschischang, Frey, and Loeliger] F. R. Kschischang, B. J. Frey, and H. . Loeliger. 2001. Factor graphs and the sum-product algorithm. *IEEE Transactions on Information Theory*, 47(2):498–519. ISSN 0018-9448.
- [Kuleshov and Ermon(2017)] Volodymyr Kuleshov and Stefano Ermon. 2017. Neural variational inference and learning in undirected graphical models. *CoRR*, abs/1711.02679. URL <http://arxiv.org/abs/1711.02679>.
- [Kullback and Leibler(1951)] S. Kullback and R. A. Leibler. 1951. On information and sufficiency. *Ann. Math. Statist.*, 22(1):79–86. URL <https://doi.org/10.1214/aoms/1177729694>.
- [Kullback(1959)] Solomon Kullback. 1959. *Information Theory and Statistics*. Wiley, New York.
- [Lazaro-Gredilla *et al.*(2019)] Lazaro-Gredilla, Lehrach, and George] Miguel Lazaro-Gredilla, Wolfgang Lehrach, and Dileep George. 2019. Learning undirected models via query training.
- [Lin *et al.*(2016)] Lin, Shen, van den Hengel, and Reid] Guosheng Lin, Chunhua Shen, Anton van den Hengel, and Ian Reid. 2016. Efficient piecewise training of deep structured models for semantic segmentation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

- [Lu *et al.*(2018) Lu, Liu, and Huang] You Lu, Zhiyuan Liu, and Bert Huang. 2018. Block belief propagation for parameter learning in markov random fields. *CoRR*, abs/1811.04064. URL <http://arxiv.org/abs/1811.04064>.
- [Montavon *et al.*(2016) Montavon, Müller, and Cuturi] Grégoire Montavon, Klaus-Robert Müller, and Marco Cuturi. 2016. Wasserstein training of restricted boltzmann machines. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 3718–3726. Curran Associates, Inc.
- [Neath(2012)] Ronald C. Neath. 2012. On convergence properties of the monte carlo em algorithm.
- [Oppor and Saad(2001)] M. Oppor and D. Saad. 2001. *Advanced Mean Field Methods: Theory and Practice*. Neural information processing series. MIT Press. ISBN 9780262150545.
- [Qu *et al.*(2019) Qu, Bengio, and Tang] Meng Qu, Yoshua Bengio, and Jian Tang. 2019. Gmn: Graph markov neural networks. In *International Conference on Machine Learning*, pages 5241–5250.
- [Salimans *et al.*(2016) Salimans, Goodfellow, Zaremba, *et al.*] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, *et al.* 2016. Improved techniques for training gans. In *Advances in Neural Information Processing Systems 29*, pages 2234–2242. Curran Associates, Inc.
- [Santambrogio(2015)] F. Santambrogio. 2015. *Optimal Transport for Applied Mathematicians: Calculus of Variations, PDEs, and Modeling*. Progress in Non-linear Differential Equations and Their Applications. Springer International Publishing. ISBN 9783319208282. URL <https://books.google.se/books?id=UOHHCGAAQBAJ>.
- [Srikumar *et al.*(2012) Srikumar, Kundu, and Roth] Vivek Srikumar, Gourab Kundu, and Dan Roth. 2012. On amortizing inference cost for structured prediction. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 1114–1124, Jeju Island, Korea, July 2012. Association for Computational Linguistics.
- [Sutton and McCallum(2012)] Charles A. Sutton and Andrew McCallum. 2012. Piecewise training for undirected models. *CoRR*, abs/1207.1409. URL <http://arxiv.org/abs/1207.1409>.
- [Tolstikhin *et al.*(2017) Tolstikhin, Bousquet, Gelly, and Schoelkopf] I. Tolstikhin, O. Bousquet, S. Gelly, and B. Schoelkopf. 2017. Wasserstein Auto-Encoders. *ArXiv e-prints*.

- [Wainwright and Jordan(2008)] M. J. Wainwright and M. I. Jordan. 2008. *Graphical Models, Exponential Families, and Variational Inference*. now. URL <https://ieeexplore.ieee.org/document/8187302>.
- [Wainwright(2006)] Martin J. Wainwright. 2006. Estimating the "wrong" graphical model: Benefits in the computation-limited setting. *J. Mach. Learn. Res.*, 7: 1829–1859. URL <http://jmlr.org/papers/v7/wainwright06a.html>.
- [Wiseman and Kim(2019)] Sam Wiseman and Yoon Kim. 2019. Amortized bethe free energy minimization for learning mrfs. In *Advances in Neural Information Processing Systems 32*, pages 15520–15531. Curran Associates, Inc.
- [Yedidia *et al.*(2003)] Yedidia, Freeman, and Weiss] Jonathan Yedidia, William Freeman, and Yair Weiss. 2003. *Understanding belief propagation and its generalizations*, volume 8, pages 239–269. ISBN 1558608117.