

Factor Graph Tutorial

HU, Pili

March 2, 2012

Abstract

Contents

1	A Motivating Example	3
1.1	Marginalization Problem	3
1.2	Inference Problem	3
1.3	Inference with Four Binary Variables	4
1.3.1	Search Max on Joint Distribution Directly	4
1.3.2	Search Max Intelligently	5
1.3.3	Observation	6
1.4	Inference with One Observed Variable	7
1.4.1	Naive	7
1.4.2	Intelligent	7
1.4.3	Observation	8
2	Factor Graph Specification	9
3	Factor Graph Transformation	9
4	Conversion from Other Model	9
5	Discussions	9

1 A Motivating Example

1.1 Marginalization Problem

Consider a joint probability distribution:

$$p(\vec{x}) \tag{1}$$

where $\vec{x} = \{x_1, x_2, \dots, x_n\}$.

Marginalization operator:

$$p_1(x_1) = \sum_{x_2} \dots \sum_{x_{n-1}} \sum_{x_n} p(\vec{x}) \tag{2}$$

Assuming discrete variable. For continuous ones, substitute sum with integral accordingly.

For simplicity of notation, introduce the shorthand "summary" notation:

$$p_1(x_1) = \sum_{\sim\{x_1\}} p(\vec{x}) \tag{3}$$

$$= \sum_{\{x_2, x_3, \dots, x_n\}} p(\vec{x}) \tag{4}$$

$$= \sum_{x_2} \dots \sum_{x_{n-1}} \sum_{x_n} p(\vec{x}) \tag{5}$$

The marginalization problem is defined as: Given $p(\vec{x})$, find $p_i(x_i)$. This problem can be generalized as summary for more than one variables.

1.2 Inference Problem

The inference problem is defined as: Given $p(\vec{x}, \vec{y})$, and the observed value of \vec{y} , say $\hat{\vec{y}} = \{\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n\}$, find the most probable configuration of \vec{x} :

$$\vec{x}^* = \arg \max_{\vec{x}} \{p(\vec{x}|\hat{\vec{y}})\} \tag{6}$$

The conditional probability can be rewritten as:

$$p(\vec{x}|\hat{\vec{y}}) = \frac{p(\vec{x}, \hat{\vec{y}})}{p(\hat{\vec{y}})} \tag{7}$$

$$p(\vec{x}|\hat{\vec{y}}) \propto p(\vec{x}, \hat{\vec{y}}) \tag{8}$$

Thus eqn(6) can be rewritten as:

$$\vec{x}^* = \arg \max_{\vec{x}} \{p(\vec{x}, \hat{\vec{y}})\} \tag{9}$$

We'll bridge the gap between the inference problem and marginalization problem defined above later. Now, we start with a toy example.

1.3 Inference with Four Binary Variables

Assume we have a joint distribution $p(a, b, c, d)$. Without any knowledge of the internal structure of the distribution, we can always write it as:

$$p(a, b, c, d) = p(a)p(b|a)p(c|a, b)p(d|a, b, c) \quad (10)$$

Now assume the distribution can be factorized in the following way:

$$p(a, b, c, d) = p(a)p(b)p(c|b)p(d|c) \quad (11)$$

We'll compare two ways of computing

$$\max_{abcd} p(abcd)$$

using the following data:

$$p(a) = \begin{bmatrix} 0.1 & 0.9 \end{bmatrix} \quad (12)$$

$$p(b) = \begin{bmatrix} 0.2 & 0.8 \end{bmatrix} \quad (13)$$

$$p(c|b) = \begin{bmatrix} 0.4 & 0.6 \\ 0.6 & 0.4 \end{bmatrix} \quad (14)$$

$$p(d|c) = \begin{bmatrix} 0.3 & 0.7 \\ 0.8 & 0.2 \end{bmatrix} \quad (15)$$

1.3.1 Search Max on Joint Distribution Directly

First, we pretend there is no structure information available. Thus a naive way to compute the max is to evaluate the joint distribution everywhere on the complete alphabets of variables, and then get maximum by comparison.

$$p(abcd) = \begin{bmatrix} \text{abcd} & \text{Probability} \\ 0000 & 0.0024 \\ 0001 & 0.0056 \\ 0010 & 0.0096 \\ 0011 & 0.0024 \\ 0100 & 0.0144 \\ 0101 & 0.0336 \\ 0110 & 0.0256 \\ 0111 & 0.0064 \\ 1000 & 0.0216 \\ 1001 & 0.0504 \\ 1010 & 0.0864 \\ 1011 & 0.0216 \\ 1100 & 0.1296 \\ 1101 & 0.3024 \\ 1110 & 0.2304 \\ 1111 & 0.0576 \end{bmatrix} \quad (16)$$

$$p(abcd^*) = \max_{abcd} p(abcd) \quad (17)$$

$$= p(1101) \quad (18)$$

$$= 0.3024 \quad (19)$$

Corresponding computation complexity:

- Function evaluation: $16 \times 4 = 64$
- Product: $16 \times 3 = 48$
- Comparison(for max operator): 15

1.3.2 Search Max Intelligently

Indeed, eqn(11) conveys useful information by the factorization of the joint probability.

Let's expand the maximization $p(a, b, c, d)$

$$\max_{abcd} \{p(abcd)\} = \max_{abcd} \{p(a)p(b)p(c|b)p(d|c)\} \quad (20)$$

$$= \max_a \{p(a)\} \max_{bcd} \{p(b)p(c|b)p(d|c)\} \quad (21)$$

$$= \max_a \{p(a)\} \max_d \{\max_{bc} \{p(b)p(c|b)p(d|c)\}\} \quad (22)$$

$$= \max_a \{p(a)\} \max_d \{\max_c \{\max_b \{p(b)p(c|b)\}p(d|c)\}\} \quad (23)$$

$$\max_a \{p(a)\} = \max_a f_a(a) = 0.9 \quad (24)$$

$$\max_b \{p(b)p(c|b)\} = \max_b f_{bc}(bc) \quad (25)$$

$$= \max_b \left[\begin{array}{c|c} bc & \text{Probability} \\ \hline 00 & 0.08 \\ 01 & 0.12 \\ 10 & 0.48(*) \\ 11 & 0.32(*) \end{array} \right] \quad (26)$$

$$= \left[\begin{array}{c|c} c & \text{Probability} \\ \hline 0 & 0.48 \\ 1 & 0.32 \end{array} \right] \quad (27)$$

Denote $\max_b \{p(b)p(c|b)\}$ by $\mu_{bc}(c)$.

$$\max_c \{\mu_{bc}(c)p(d|c)\} = \max_c f_{cd}(cd) \quad (28)$$

$$= \max_c \left[\begin{array}{c|c} cd & \text{Probability} \\ \hline 00 & 0.144 \\ 01 & 0.336^{(*)} \\ 10 & 0.256^{(*)} \\ 11 & 0.064 \end{array} \right] \quad (29)$$

$$= \left[\begin{array}{c|c} d & \text{Probability} \\ \hline 0 & 0.256 \\ 1 & 0.336 \end{array} \right] \quad (30)$$

Denote $\max_c \{\mu_{bc}(c)p(d|c)\}$ by $\mu_{cd}(d)$.

$$\max_d \{\max_c \{\max_b \{p(b)p(c|b)\}p(d|c)\}\} \quad (31)$$

$$= \max_d \{\mu_{cd}(d)\} \quad (32)$$

$$= 0.336 \quad (33)$$

Thus we get final result:

$$\max_{abcd} \{p(abcd)\} = \max_a \{p(a)\} \times \max_d \{\mu_{cd}(d)\} \quad (34)$$

$$= 0.3024 \quad (35)$$

Again, we calculate the computation complexity:

- Function evaluation:

$$2 + 4 \times 2 + 4 \times 2 + 2 = 20$$

- Product:

$$0 + 4 \times 1 + 4 \times 1 + 0 + 1 = 9$$

- Comparison(for max operator):

$$1 + 2 \times 1 + 2 \times 1 + 1 = 6$$

1.3.3 Observation

We compare the complexity of two methods in table(1).

The observations:

- By properly using the structure of joint distribution, it's possible to reduce computation complexity.

Table 1: Comparison Between Two Methods

Items	Naive	Intelligent
Function	64	20
Product	48	9
Comparison	15	6

- The trick to reduce complexity in the second method is: "product" is distributive through "max". Thus we can separate some variables when evaluating the maximum of others. We'll address this issue in details later.
- How to reveal and utilize the structure in a systematic way is still a problem.

1.4 Inference with One Observed Variable

Here's another example. Assume c is observed to be 1 in the last example. What's the new most probable configuration of other variables?

1.4.1 Naive

As before, simply restrict the evaluation of functions only on points where $c = 1$.

$$p(abcd) = \begin{array}{c|c} \text{abcd} & \text{Probability} \\ \hline 0010 & 0.0096 \\ 0011 & 0.0024 \\ 0110 & 0.0256 \\ 0111 & 0.0064 \\ 1010 & 0.0864 \\ 1011 & 0.0216 \\ 1110 & 0.2304 \\ 1111 & 0.0576 \end{array} \quad (36)$$

The most probable configuration of the four variables is 1110, and corresponding probability is 0.2304 (the joint probability, not the probability of $a = 1, b = 1, d = 0$ conditioned $c = 1$).

1.4.2 Intelligent

With the observation of $c = 1$, the joint distribution can be decomposed as:

$$\max_{abd} \{p(ab1d)\} = \max_{abd} \{p(a)p(b)p(c = 1|b)p(d|c = 1)\} \quad (37)$$

$$= \max_a \{p(a)\} \max_b \{p(b)p(c = 1|b)\} \max_d \{p(d|c = 1)\} \quad (38)$$

$$\max_a \{p(a)\} = \max_a f_a(a) = 0.9 \quad (39)$$

$$\max_b \{p(b)p(c=1|b)\} = \max_b f_{bc}(bc, c=1) \quad (40)$$

$$= \max_b \left[\begin{array}{c|c} bc & \text{Probability} \\ \hline 01 & 0.12 \\ 11 & 0.32 \end{array} \right] \quad (41)$$

$$= 0.32 \quad (42)$$

$$\max_d \{p(d|c=1)\} = \max_d \left[\begin{array}{c|c} cd & \text{Probability} \\ \hline 10 & 0.8 \\ 11 & 0.2 \end{array} \right] \quad (43)$$

$$= 0.8 \quad (44)$$

Thus the final maximum probability is given by:

$$\max_a \{p(a)\} \max_b \{p(b)p(c=1|b)\} \max_d \{p(d|c=1)\} \quad (45)$$

$$= 0.9 * 0.32 * 0.8 \quad (46)$$

$$= 0.2304 \quad (47)$$

1.4.3 Observation

Now that we validated the correctness of our intelligent method, we again compare the complexity as is in table(2).

Table 2: Comparison Between Two Methods

Items	Naive	Intelligent
Function	32	8
Product	24	4
Comparison	7	3

Besides previous observations on the value of "structure", we highlight one more thing:

- When the variable c is observed, the joint distribution function can be further decomposed! That is, in previous example, there is a relationship between b and d , so we evaluate max operator in the order of b , then c , then d . However, with the observation of c , the sub functions involving b and d are fully decoupled, this further reduces the complexity.
- This observation is indeed the notion of conditional independence, as is one major concern in some graphical models like MRF and BN.

2 Factor Graph Specification

3 Factor Graph Transformation

4 Conversion from Other Model

5 Discussions

Since the development of factor graph is boosted in the past decade, different authors come up with different description of similar problems. Not to distinguish right from wrong, I just regard those stuffs out there as inconsistent. My opinion on some parts of past literature:

- In Bishop's book[1], chapter 8.4.5, P411. The example is not good. Actually, when talking about that probability maximization problem, we should know "product" corresponds to product operator, and "sum" corresponds to max operator. In this case, the marginalization operation for a single variable is indeed the maximization for each instance of that variable. Using local marginalized function(max), we can certainly get the global probability maximization point considering all variables.
- As for Dynamic Factor Graph, the author of this paper do not advocate the abuse of this term like an extension of factor graph. FG itself is able to model system dynamics, as we've already seen in those examples above. Other authors may use the term DFG [5] [4] , but their DFG is application specific. Those graphs are essentially FG. Not until we examine the physical meaning of some factor nodes do we realize their "dynamic" property.

References

- [1] C.M Bishop. *Pattern recognition and machine learning*, volume 4. springer New York, 2006.
- [2] B.J. Frey. Extending factor graphs so as to unify directed and undirected graphical models. In *Proc. 19th Conf. Uncertainty in Artificial Intelligence*, pages 257–264, 2003.
- [3] F.R. Kschischang, B.J. Frey, and H.A. Loeliger. Factor graphs and the sum-product algorithm. *Information Theory, IEEE Transactions on*, 47(2):498–519, 2001.
- [4] P. Mirowski and Y. LeCun. Dynamic factor graphs for time series modeling. *Machine Learning and Knowledge Discovery in Databases*, v:128–143, 2009.

- [5] C. Wang, J. Tang, J. Sun, and J. Han. Dynamic social influence analysis through time-dependent factor graphs. *ASONAM*, v:p, 2011.
- [6] Moral Graph, Wikipedia, http://en.wikipedia.org/wiki/Moral_graph
- [7] Markov Blanket, Wikipedia, http://en.wikipedia.org/wiki/Markov_blanket