

Factor Graph Tutorial

HU, Pili

March 2, 2012

Abstract

Factor graph is one general problem solving technique (personally, I prefer to call it a technique rather than a model). It can solve marginalization problem very efficiently. Our traditional model, which involve probability marginalization, like Hidden Markov Model, Bayesian Network, Markov Random Field, etc, can all be translated into a factor graph representation. This article first constructs a toy motivating example to show the benefit of using factor graph. With the observation from the examples, factor graph is deduced. Notations, definitions, and descriptions mainly follow the original works, [4] [3]. Then we present two selected applications of FG, one for modeling underlying system dynamics[5] , and the other for modeling influence in social network[6]. We end this article by relating factor graphs to other existing graphical models and algorithms, together with some discussion.

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	9
2.1	Specification	9
2.2	Transformation	9
3	Selected Applications	10
3.1	Model System Dynamics	10
3.1.1	Quick Note of the Main Idea	10
3.1.2	Comments	12
3.2	Solving Marginalization Problem in Social Network	13
3.2.1	Quick Note of the Main Idea	13
3.2.2	Comments	15
3.3	Short Summary on the Application of FG	16
4	Other Taste of Modeling	16
5	Related Models/Algorithms	16
6	Discussions	18

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.

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

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

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

$$= 0.3024 \quad (19)$$

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

2.1 Specification

2.2 Transformation

Transformation of loopy factor graph can result in tree structures. Thus ordinary sum-product algorithm can be used to obtain an exact solution.

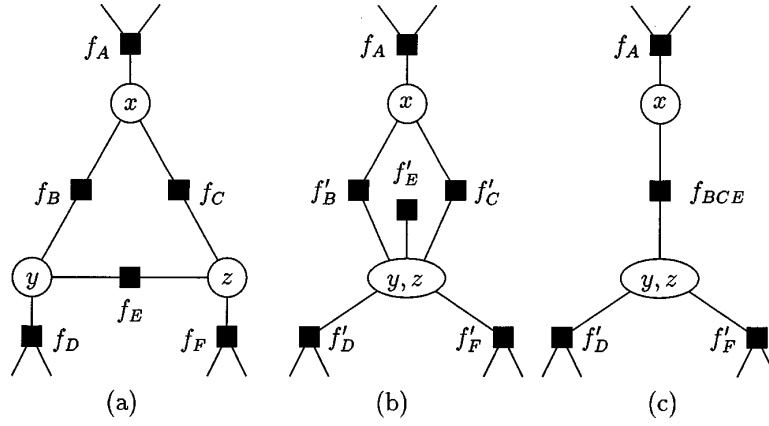


Figure 1: Transformation: Clustering[4]

Effect of clustering:

- Cluster variable nodes: domain enlarged; no complexity increasing in local functions; increase complexity in sum-product algorithm (function evaluation).
- Cluster function nodes: do not increase complexity of variables; increase complexity in sum-product algorithm (sizes of messages are increased).

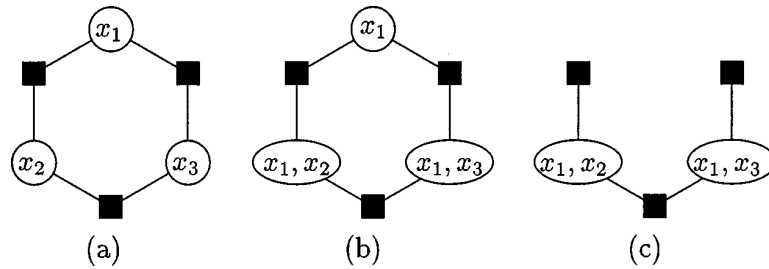


Figure 2: Transformation: Stretching[4]

Effect of stretching:

- Simply stretching variables adds positional variables to some local functions. Local function complexity is not changed. Variable alphabets are enlarged.
- After stretching, redundant links or even redundant variables can be removed. By systematic stretching, all cycles can be eliminated.

3 Selected Applications

3.1 Model System Dynamics

This section follows the paper:

P. Mirowski and Y. LeCun. Dynamic factor graphs for time series modeling. *Machine Learning and Knowledge Discovery in Databases*, v:128–143, 2009.

Declaration: figures in this section are borrowed from the original paper[5]. We omit further citation for simplicity.

3.1.1 Quick Note of the Main Idea

Problem settings:

- System is composed of state variables and observation variables.
- Observation function maps state to observation.
- Dynamic function governs state transitions.
- Aim at modeling high dimensional underlying dynamic, probably non-linear, but deterministic.

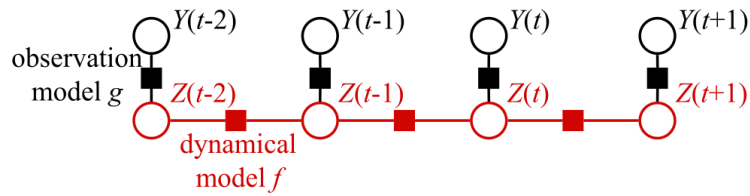


Figure 3: Single State Dependency(HMM)

When current state is only dependent on previous one state, it is like an HMM model.

The state transition model can be more general that it may depend on several previous states, or even observations.

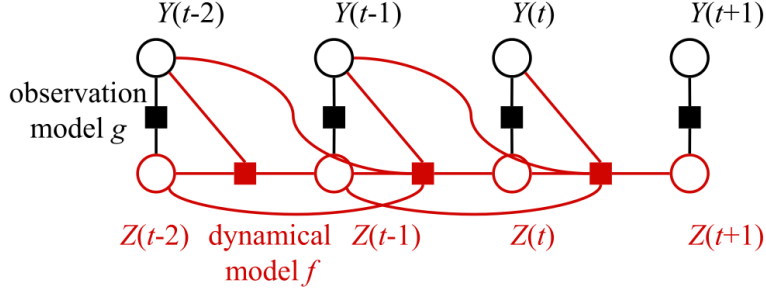


Figure 4: Multiple States and Observation Dependency

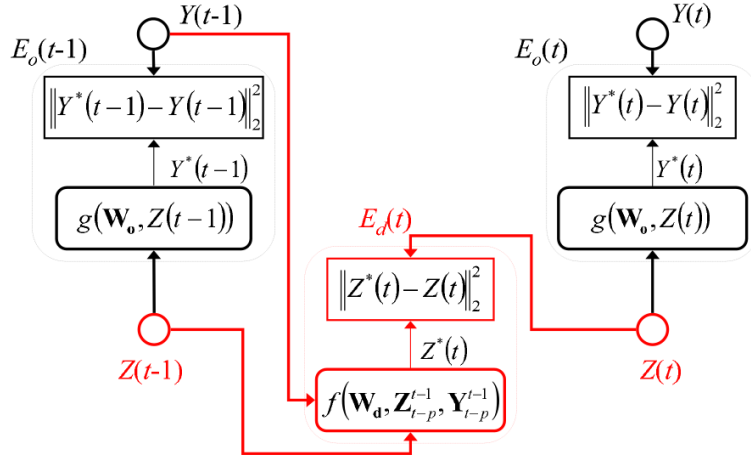


Figure 5: System Details

Formulation of parameterized FG:

$$L(\mathbf{W}, Y, Z) = E(\mathbf{W}, Y) + R_z(Z) + R(\mathbf{W}) \quad (48)$$

$$\tilde{Z} = \arg \max_Z L(\tilde{\mathbf{W}}, Y, Z) \quad (49)$$

$$\tilde{\mathbf{W}} = \arg \max_{\mathbf{W}} L(\mathbf{W}, Y, \tilde{Z}) \quad (50)$$

Since the paper operates in the negative log domain, they want to minimize L . The more likely some configuration is, the energy E is lower. As is depicted in fig(), energy function is the square error.

When parameters are fixed, eqn(49) acts as what ordinary factor graph does, namely, given all explicit function definitions, compute the marginalization problem. In this application, parameters \mathbf{W} is not determined. An usual way to tackle with this problem is by Expectation-Maximization schema (sometimes called Generalized Expectation Maximization algorithm). eqn(49) resembles the E-step, while eqn(50) resembles the M-step.

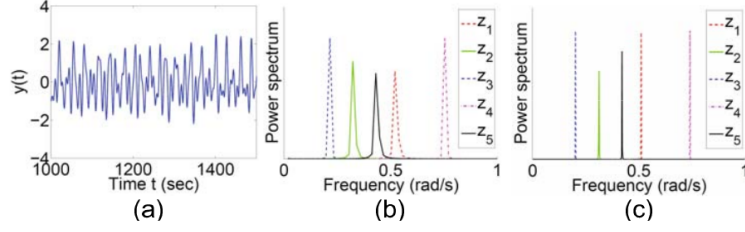


Figure 6: Training and Evaluation on Superimposed Sine Waves

To decouple 5 superimposed sine waves, the authors choose the architecture:

- 5 dimensional state variable.
- Dynamic function(f): 5 independent FIR filter of order 25.

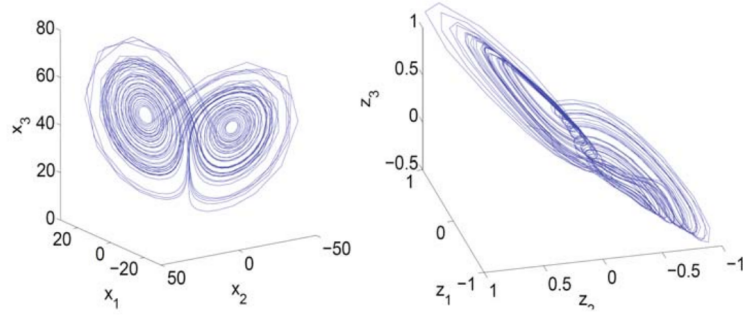


Figure 7: Training and Evaluation on Lorenz Curve

To capture the observation of 1D variable whose underlying dynamic is 3D-Lorenz curve, the authors choose the architecture:

- 3 dimensional state variable.
- Dynamic function(f): 3 layered convolutional network.

Besides the two evaluation mentioned above, the original paper contains more. Interested readers can refer to [5] for more information.

3.1.2 Comments

- DFG(in this paper) is basically FG. "dynamic" only address the application aspects, rather than claiming an extension of FG.
- The framework of FG is too general. Designing observation function and state transition function, choice of regularizer, setting coefficients are really tricky. There doesn't seem to be a rule of thumb.

- The introduction of unknown parameters and regularizers make the formulation deviate from ordinary factor graph. Before final marginalization, there is a training stage to obtain optimal parameters.

3.2 Solving Marginalization Problem in Social Network

This section follows the paper:

C. Wang, J. Tang, J. Sun, and J. Han. Dynamic social influence analysis through time-dependent factor graphs. *ASONAM*, v:p, 2011.

Declaration: figures in this section are borrowed from the original paper [6]. We omit further citation for simplicity.

3.2.1 Quick Note of the Main Idea

Problem settings:

- Graph $G^t = \langle V^t, E^t \rangle$. V corresponds to people in a certain social network. E corresponds activity relationships. There is a real valued weight function w defined on E .
- Graph G^t captures time varying relationships.
- Definition of w is problem(application) specific.
- Want to find the pairwise influence μ_{ij} . Denote the node to influence node i by y_i , a random variable. μ_{ij} is $P(y_i = j)$.

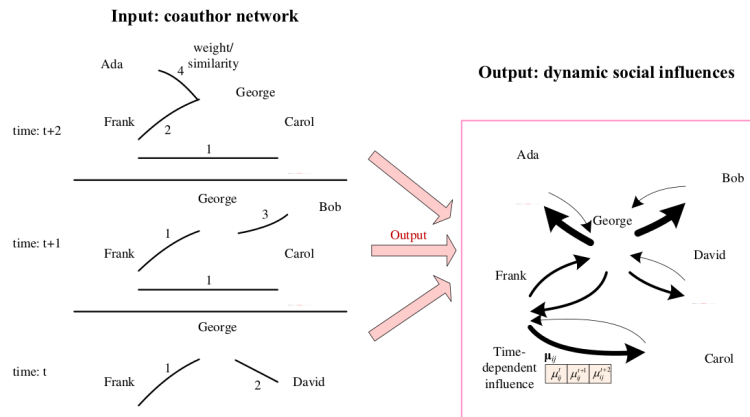


Figure 8: Problem Definition

When defined on the static version, one μ_{ij} is output for each i, j . When defined on the dynamic version, a vectore $(\mu_{ij}^{(t)})$ is output.

Two definitions to facilitate further description:

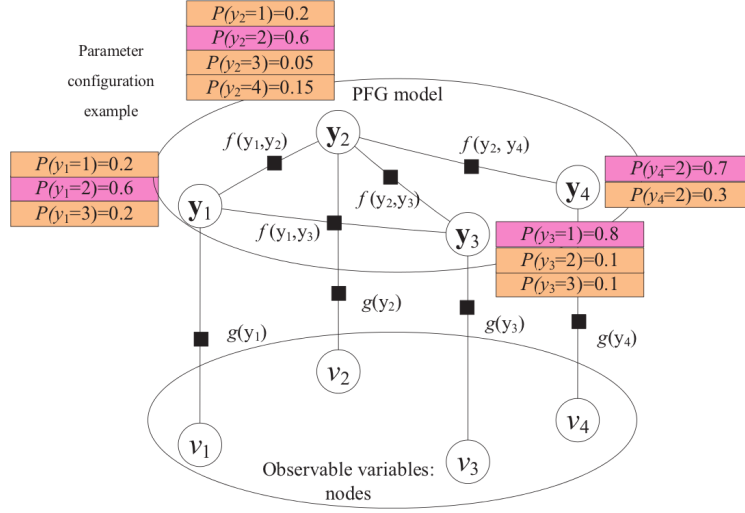


Figure 9: Visualization of the Factor Graph

- $NB(i)$: the neighbourhood of i . (More precisely, v_i , this is defined on the original graph)
- $SC(i) = NB(i) \cup i$.

More precisely, y_i should be linked to $v_j, \forall j \in NB(i)$ in the FG depicted in fig(). This can be justified through the definition of node factor function.

$$g(y_i|V) = \begin{cases} \frac{w_{iy_i}}{\sum_{j \in NB(i)} (w_{ij} + w_{ji})} & y_i \neq i \\ \frac{\sum_{j \in NB(i)} w_{ji}}{\sum_{j \in NB(i)} (w_{ij} + w_{ji})} & y_i = i \end{cases}$$

Figure 10: Choice of Node Factor Function

$$f(y_i, y_j|V) = \begin{cases} \frac{u}{|SC(i) \cap SC(j)|} & y_i = y_j \\ \frac{1-u}{|SC(i)| |SC(j)| - |SC(i) \cap SC(j)|} & y_i \neq y_j \end{cases}$$

Figure 11: Choice of Link Factor Function

The definition of node factor function and link factor function, together with the observation of G removes complex links from y_i to V . The resultant FG can be viewed as one single layer with y_i vertices as variable nodes. Function nodes $g_i(y_i)$ are completely defined with the observation of G .

Note that there is one parameter u in the choice of link factor function. If not for this paramter, the FG in this paper was able to marginalize each y_i directly and give the influence using $\mu_{ij} = P(y_i = j)$. Again, the parameter estimation can be done using EM schema.

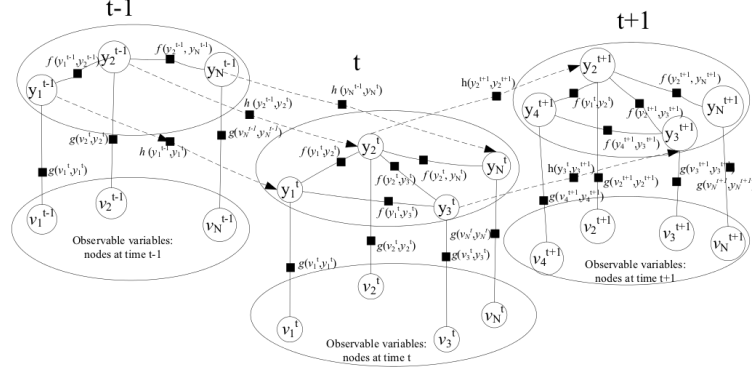


Figure 12: Factor Graph with Time Dependency

In order to capture the time varying property of social influence, the authors augmented original static FG to be time-dependent FG. Bridging function nodes are added between two time consecutive sub graphs. For detailed discussion, please refer to original paper.

Since the current section aims at modeling, details of algorithms are omitted. However, it's worth referring to the original paper for a specialized message passing algorithm named SIP(Social Influence Propagation) by the authors. That algorithm works more efficiently than general junction tree algorithms.

3.2.2 Comments

- Modeling an FG is rather easy. However, coming up with right choice of functions and paramters are difficult.
- In this paper, the authors used heuristics to generate w_{ij} , which is a very important raw data. After the calculation of w_{ij} , the FG and so does the marginalizations are fully defined.
- The original intention(we've already shown that at the beginning of this article) of FG is to reduce complexity by utilizing graphical structures. However, the compuation in practice is still a concern. How to carefully design the algorithm to perform marginalization under the message passing framework according to physical meanings will be an interesting and channelling problem.

3.3 Short Summary on the Application of FG

- Framework and general algorithm is mature both in theory and in practice.
- Original FG can be used to solve marginalization problem directly. At this point, FG can be think of as a problem solver rather than a model. For this flavour, please refer to Bishop's book [1], Chapter 8. The author begins with BN and MRF which model directed and undirected independencies. In order to solve the inference problem, the author propose to convert BN and MRF to FG.
- According to different backgrounds, different authors come up with different Paramterized Factor Graph. Most of the works include two efforts:
 - Heuristics to determine paramters.
 - Iterative methods to determine paramters. (EM schema)

Solving marginalization in FG is one sub problem then. Before giving the final marginalized answer, there is a training stage.

- Specific message passing algorithms are worth investigating. As is in [6], the authors propose SIP, specialized for their problem. The resulting algorithm is more tractable than general solvers.

4 Other Taste of Modeling

Besides modeling probability problems, people also use FG to model certain codes.

Consider a parity checking matrix given by H :

$$H = \begin{bmatrix} 1 & 1 & 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 \\ 1 & 0 & 1 & 1 & 0 & 0 \end{bmatrix} \quad (51)$$

In the corresponding factor graph, each factor node corresponds to one row and each variable node corresponds to one column. The type of the functions are indicators testing the predicate: $[x_{i_1} \oplus x_{i_2} \oplus x_{i_3} = 0]$. Note that the logic " \cap " operator can be regarded as "product". This kind of modeling is called behavioural modeling in [4].

5 Related Models/Algorithms

Models:

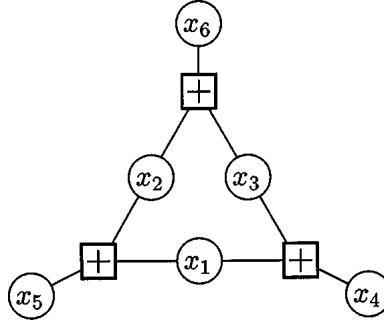


Figure 13: A Parity Code Modeled by FG[4]

- Bayesian Network. Directed Graphical Model.
- Markov Random Field. Undirected Graphical Model.
- Chain Graph. Directed + Undirected.
- Factor Graph.

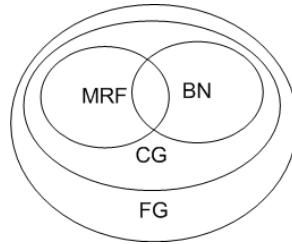


Figure 14: Venn Graph of Several Graphical Models

Justification of modeling ability:

- $BN \rightarrow MRF$. The moralization process hides certain structure previously available in BN.
- MRF. An example of square located 4 nodes. If two diagonal variables are observed, the other two variables become independent. BN can not express such information.
- CG. Since both directed and undirected, it forms the superset of the union of BN and MRF.
- FG. Converting CG to FG reveals more details, namely, how the joint distribution is factorized(usually, one explicit normalization function exhibits). BTW, FG can model not only probabilistic problems, but also other problems like system dynamics and coding. It is super general in the sense.

Algorithms:

- Junction tree.
- Bayesian ball.
- (Loopy)Belief propagation.
- Sum-product.

6 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 [6] [5] , 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.

Acknowledgements

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] B.J. Frey, F.R. Kschischang, H.A. Loeliger, and N. Wiberg. Factor graphs and algorithms. In *PROCEEDINGS OF THE ANNUAL ALLERTON CONFERENCE ON COMMUNICATION CONTROL AND COMPUTING*, volume 35, pages 666–680. UNIVERSITY OF ILLINOIS, 1997.
- [4] 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.
- [5] P. Mirowski and Y. LeCun. Dynamic factor graphs for time series modeling. *Machine Learning and Knowledge Discovery in Databases*, v:128–143, 2009.
- [6] C. Wang, J. Tang, J. Sun, and J. Han. Dynamic social influence analysis through time-dependent factor graphs. *ASONAM*, v:p, 2011.
- [7] Moral Graph, Wikipedia, http://en.wikipedia.org/wiki/Moral_graph
- [8] Markov Blanket, Wikipedia, http://en.wikipedia.org/wiki/Markov_blanket