

## Efficient computation of maximum likelihood estimators of hierarchical subspace models

Hisayuki Hara \*

Faculty of Economics, Niigata University, Niigata, Japan  
hara@econ.niigata-u.ac.jp

### Abstract

In this paper we discuss efficient computation of the maximum likelihood estimator (MLE) for general log linear models for contingency tables. In general an iterative algorithm, such as the iterative proportional fitting (IPF), is required for computing the MLE of a log linear model. For a hierarchical model, it is well known that the updating rule for IPF is localized according to the separation or decomposition of the graph describing conditional independence structure among variables of the model and then the computational cost is reduced (e.g. [1]). Here we extend the algorithm for a hierarchical model to the hierarchical subspace model (HSM) introduced by [2] and propose an efficient algorithm of updating rule of IPF for HSM.

**Keywords** : chordal extension, graphical model, hierarchical model, information propagation, iterative proportional fitting

## 1 Introduction

Consider an  $m$ -way contingency table  $\{x(\mathbf{i}), \mathbf{i} \in \mathcal{I}\}$  with the set of variables  $[m] = \{1, \dots, m\}$ , where  $\mathbf{i}$  denotes each cell and  $\mathcal{I}$  denotes the set of cells.  $x(\mathbf{i})$  is a frequency for a cell  $\mathbf{i}$ . For a subset  $D \subset [m]$ , denote by  $\mathbf{i}_D$  a marginal cell for  $D$ . Let  $\mathcal{I}_D$  be the set of marginal cells for  $D$ . Let  $\Delta$  be a simplicial complex with the set of vertices  $[m]$ . Then hierarchical model defined by  $\Delta$  is expressed by

$$\log p(\mathbf{i}) = \sum_{D \in \Delta} \mu_D(\mathbf{i}_D). \tag{1}$$

Models (1) with linear restrictions among natural parameters is extensively used in practical data analysis (e.g. [3, 4]). The canonical form of such a model is written by

$$\log p(\mathbf{i}) = \sum_{D \in \tilde{\Delta}} \phi_D(\mathbf{i}_D) \beta_D, \tag{2}$$

where  $\phi_D(\mathbf{i}_D) = \{\phi_D^k(\mathbf{i}_D)\}_{k=1, \dots, K_D}$  is  $1 \times K_D$  vector of known functions of  $\mathbf{i}_D$  and  $\beta_D$  is  $K_D \times 1$  parameter vector. Since the model (2) is the models (1) with linear constraint, these two models have the same interaction terms corresponding to  $\Delta$ . In canonical form (2), however, we note that  $\tilde{\Delta}$  is also a simplicial complex but not equal to  $\Delta$  in general (e.g. [2]). The model (2) is still an exponential family and therefore a linear subspace of hierarchical model (1). The model (2) is also called hierarchical subspace model (HSM, [2]).

In this paper we discuss an efficient computation of maximum likelihood estimator (MLE) of HSM (2). Define

$$t_D := \sum_{\mathbf{i}_D \in \mathcal{I}_D} \phi_D(\mathbf{i}_D) x(\mathbf{i}_D).$$

Then the sufficient statistic  $t$  for HSM (2) is

$$t = \{t_D \mid D \in \tilde{\Delta}\}.$$

Then the likelihood equations are written by

$$t_D = E[t_D] = n \sum_{\mathbf{i}_D \in \mathcal{I}_D} \phi_D(\mathbf{i}_D) p(\mathbf{i}_D), \quad D \in \tilde{\Delta}.$$

In general MLE of HSM (2) is not explicitly obtained and iterative computation is required to compute MLE. Iterative proportional fitting algorithm (IPF) is one of the popular algorithms for computing MLE. IPF for HSM is described as follows.

**Algorithm 1** (IPS for HSM).

**Step 0**  $t \leftarrow 0, p^{(t)}(\mathbf{i}) \leftarrow 1/n$

**Step 1** For all  $D \in \tilde{\Delta}$  and  $k = 1, \dots, K_D$ , update  $p^{(t)}(\mathbf{i})$  as

$$p^{(t+1)}(\mathbf{i}) \leftarrow \frac{\sum_{\mathbf{i}_D \in \mathcal{I}_D} \phi_D^k(\mathbf{i}_D) x(\mathbf{i}_D)}{n \sum_{\mathbf{i}_D \in \mathcal{I}_D} \phi_D^k(\mathbf{i}_D) p^{(t)}(\mathbf{i}_D)} \cdot p^{(t)}(\mathbf{i}), \quad \mathbf{i} \in \mathcal{I}. \quad (3)$$

**Step 2** If  $p^{(t+1)}(\mathbf{i})$  converges for all  $\mathbf{i}$ , output  $p^{(t+1)}(\mathbf{i})$ . If not, go back to Step 1.

It is well known that the convergence of this algorithm to MLE is guaranteed (e.g. [5]). Step 1 requires  $O(|\mathcal{I}|)$  times updates. For a hierarchical model (1), however, it is known that the updating rule for IPF is localized according to the separation or decomposition of the graph describing conditional independence structure among variables of the model and then the computational cost can be reduced (e.g. [1]). The main purpose of this paper is to extend the argument in hierarchical model to HSM. As discussed in [2], decomposition of an HSM does not always correspond to that of its conditional independence graph  $G_\Delta$  and hence the algorithm of [1] cannot be applied directly to HSM. In this paper we show that by defining the decomposition of a HSM properly, we can apply the algorithm of [1] also to HSM.

## 2 Main Results

### 2.1 Decomposition of HSM

As discussed the previous section, the canonical form of HSM (2) is defined by a simplicial complex  $\tilde{\Delta}$ . In general  $\tilde{\Delta}$  is not equal to  $\Delta$ . As mentioned in the previous section, the model (2) and the models (1) have the same interaction terms corresponding to  $\Delta$ . Let  $H_\Delta$  be a hypergraph with the set of vertices  $[m]$  and the set of hyperedges  $\Delta$ . Then  $H_\Delta$  also describes the conditional independence structure of both (1) and (2). Let  $\mathcal{C}$  and  $\mathcal{S}$  be the set of compact components and dividers of  $H_\Delta$  (see [1]). Let  $\nu(S)$  be the multiplicity of  $S \in \mathcal{S}$  (e.g. [5]). Then  $p(\mathbf{i})$  is expressed by

$$p(\mathbf{i}) = \frac{\prod_{C \in \mathcal{C}} p(\mathbf{i}_C)}{\prod_{S \in \mathcal{S}} p(\mathbf{i}_S)^{\nu(S)}}. \quad (4)$$

When  $\phi_D(\mathbf{i}_D) \beta_D = \mu_D(\mathbf{i}_D), \forall \mathbf{i}_D \in \mathcal{I}_D$ , in (2),  $D$  is called saturated in (2). For any hierarchical model (1), every  $S$  is known to be saturated in (1). For HSM, however,  $S$  is

not always saturated. Furthermore, if  $S$  is not saturated in (2), marginal models  $p(\mathbf{i}_C)$  and  $p(\mathbf{i}_S)$  are no longer full exponential families and therefore they are not HSMs in general. If there is a linear constraint of parameters across more than one facet,  $\tilde{\Delta} \supset \Delta$  in the canonical form (2). This fact indicates that the decomposition of HSM does not necessarily correspond to that of  $H_\Delta$  or  $G_\Delta$ .

Let  $H_{\tilde{\Delta}}$  be a hypergraph with the set of vertices  $[m]$  and the set of hyperedges  $\tilde{\Delta}$ . Let  $\mathcal{S}^*$  be the set of dividers of  $H_{\tilde{\Delta}}$  which are saturated in (2). Let  $\mathcal{C}^*$  be the set of compact component defined by  $S^*$ . The elements of  $\mathcal{C}^*$  are also called extended compact components ([2]). Then  $p(\mathbf{i})$  satisfies

$$p(\mathbf{i}) = \frac{\prod_{C^* \in \mathcal{C}^*} p(\mathbf{i}_{C^*})}{\prod_{S^* \in \mathcal{S}^*} p(\mathbf{i}_{S^*})^{\nu(S^*)}} \tag{5}$$

and marginal models  $p(\mathbf{i}_{C^*})$  satisfy

$$p(\mathbf{i}_{C^*}) = \sum_{C^* \cap D, D \in \tilde{\Delta}} \phi_D(\mathbf{i}_D) \beta_D.$$

Hence  $p(\mathbf{i}_{C^*})$  is still an HSM ([2]). Since  $S \in \mathcal{S}^*$  is saturated in (2), the MLE of  $p(\mathbf{i}_{S^*})$  is the corresponding relative marginal frequency  $x(\mathbf{i}_{S^*})/n$ , where  $n$  is total sample size. Therefore the MLE of  $p(\mathbf{i})$  is expressed by

$$\hat{p}(\mathbf{i}) = \frac{\prod_{C^* \in \mathcal{C}^*} \hat{p}(\mathbf{i}_{C^*})}{\prod_{S^* \in \mathcal{S}^*} (x_S(\mathbf{i}_{S^*})/n)^{\nu(S^*)}}. \tag{6}$$

This fact shows that in order to compute the MLE of HSM (2) by IPF, it suffices to apply IPF to each marginal model  $p(\mathbf{i}_{C^*})$   $C^* \in \mathcal{C}^*$ . Then the computational cost for a update is reduced to  $O(\sum_{C^* \in \mathcal{C}^*} |\mathcal{I}_{C^*}|)$ .

### 2.2 Information propagation algorithm for HSM

In this section we assume that an HSM (2) satisfies  $\mathcal{S}^* = \emptyset$ . Such an HSM is called prime. For a prime HSM, we can directly apply the information propagation algorithm of [1]. Let  $\bar{G}_\Delta$  be a chordal extension of conditional independence model  $G_\Delta$ . Let  $\bar{\mathcal{C}}$  and  $\bar{\mathcal{S}}$  be the set of maximal cliques and minimal vertex separators of  $\bar{G}_\Delta$ . Let  $\bar{\mathcal{T}}_C$  be a directed clique tree of  $G_\Delta$  with a root  $\bar{C}$ . Then the algorithm is described as follows.

**Algorithm 2** (Updating rule for localized IPS at time  $t + 1$ ).

**Step 1** For  $D \in \Delta$  and  $k$ , choose any  $\bar{C} \in \bar{\mathcal{C}}$  satisfying  $D \subset \bar{C}$ .

**Step 2** Update  $p(\mathbf{i}_{\bar{C}})$ ,  $\mathbf{i} \in \mathcal{I}$  by

$$p^{(t+1)}(\mathbf{i}_{\bar{C}}) \leftarrow \frac{\sum_{\mathbf{i}_D \in I_D} \phi_D^k(\mathbf{i}_D) x(\mathbf{i}_D)}{n \sum_{\mathbf{i}_D \in I_D} \phi_D^k(\mathbf{i}_D) p^{(t)}(\mathbf{i}_D)} \cdot p^{(t)}(\mathbf{i}_{\bar{C}})$$

**Step 3** For all descendants  $\bar{C}'$  of  $\bar{C}$  on  $\bar{\mathcal{T}}_C$ , update  $p(\mathbf{i}_{\bar{C}'})$  in the following way.

- 3-1. Assume that all ancestors of  $\bar{C}'$  are already updated. From a parent  $\bar{C}''$  of  $\bar{C}'$ , receive  $p^{(t+1)}(\mathbf{i}_{\bar{S}})$ , where  $\bar{S} = \bar{C}' \cap \bar{C}''$ .

3-2. Update  $p(\mathbf{i}_{\bar{C}'})$  by

$$p^{(t+1)}(\mathbf{i}_{\bar{C}'}) \leftarrow \frac{\sum_{i_D \in I_D} p^{(t+1)}(\mathbf{i}_{\bar{S}})}{\sum_{i_D \in I_D} p^{(t)}(\mathbf{i}_{\bar{S}})} \cdot p^{(t)}(\mathbf{i}_{\bar{C}'}).$$

The computational cost of the above algorithm for updating  $p(\mathbf{i})$  is  $O(\sum_{\bar{C} \in \bar{\mathcal{C}}} |\mathcal{I}_{\bar{C}}|)$ . In general  $\sum_{\bar{C} \in \bar{\mathcal{C}}} |\mathcal{I}_{\bar{C}}| < |\mathcal{I}|$ .

**Theorem 1.** *The output of Algorithm 2 is the same as the output of the updating rule (3).*

*Proof.*  $p^{(t)}(\mathbf{i})$  is written by

$$p^{(t)}(\mathbf{i}) = \frac{\prod_{\bar{C} \in \bar{\mathcal{C}}} p^{(t)}(\mathbf{i}_{\bar{C}})}{\prod_{\bar{S} \in \bar{\mathcal{S}}} p^{(t)}(\mathbf{i}_{\bar{S}})}.$$

The updating rule for  $p^{(t)}(\mathbf{i}_{\bar{C}})$  in Step 2 is

$$p^{(t+1)}(\mathbf{i}_{\bar{C}}) = \frac{\sum_{i_D \in I_D} \phi_D^k(\mathbf{i}_D) x(\mathbf{i}_D)}{n \sum_{i_D \in I_D} \phi_D^k(\mathbf{i}_D) p^{(t)}(\mathbf{i}_D)} \cdot p^{(t)}(\mathbf{i}_{\bar{C}}),$$

The updating rule for  $p^{(t)}(\mathbf{i}_{\bar{C}'})$  in Step 3 is

$$p^{(t+1)}(\mathbf{i}_{\bar{C}'}) \leftarrow \frac{\sum_{i_D \in I_D} p^{(t+1)}(\mathbf{i}_{\bar{S}})}{\sum_{i_D \in I_D} p^{(t)}(\mathbf{i}_{\bar{S}})} \cdot p^{(t)}(\mathbf{i}_{\bar{C}'}).$$

Noting that

$$p^{(t+1)}(\mathbf{i}_S) = \frac{p^{(t+1)}(\mathbf{i}_S)}{p^{(t)}(\mathbf{i}_S)} p^{(t)}(\mathbf{i}_S),$$

we obtain

$$p^{(t+1)}(\mathbf{i}) = \frac{\sum_{i_D \in I_D} \phi_D^k(\mathbf{i}_D) x(\mathbf{i}_D)}{n \sum_{i_D \in I_D} \phi_D^k(\mathbf{i}_D) p^{(t)}(\mathbf{i}_D)} \cdot p^{(t)}(\mathbf{i}).$$

□

## References

- [1] J. H. Badsberg and F. M. Malvestuto. An implementation of the iterative proportional fitting procedure by propagation trees. *Comput. Statist. Data. Anal.*, 37:297–322, 2001.
- [2] Hisayuki Hara, Tomonari Sei, and Akimichi Takemura. Hierarchical subspace model for contingency tables. *Journal of Multivariate Analysis*, 103:19–34, 2010.
- [3] Chihiro Hirotsu. Two-way change-point model and its application. *Australian Journal of Statistics*, 39(2):205–218, 1997.
- [4] Søren Højsgaard. Split models for contingency tables. *Comput. Statist. Data. Anal.*, 42:621–645, 2003.
- [5] Steffen L. Lauritzen. *Graphical Models*. Oxford University Press, Oxford, 1996.