

# Fast and Stable Maximum Likelihood Estimation for Incomplete Multinomial Models

Chenyang Zhang, Guosheng Yin

Department of Statistics and Actuarial Science,  
The University of Hong Kong

June 13, 2019

# What is Incomplete Multinomial Model?

- A toy example: Incomplete contingency table

	Young	Middle	Senior
Female	$p_1$	$p_2$	$p_3$
Male	$p_4$	$p_5$	$p_6$

- Sample 1:

	Young	Middle	Senior
Female	21	24	18
Male	20	25	12

- Sample 2:

Female	18
Male	22

- Sample 3:

	Young	Middle	Senior
	10	20	10

- Sample 4:

	Young
Female	53
Male	47

# What is Incomplete Multinomial Model (Cont'd)

Multinomial model: the sample space  $\Omega$  is partitioned into  $K$  disjoint subspaces.  
Incomplete cases:

- (a) a subset of categories rather than a unique category is reported (partial classification).
- (b) the set of possible outcomes contains only part of all categories (truncated outcomes).

$$L(\mathbf{p}|\mathbf{a}, \mathbf{b}, \Delta) \propto \prod_{k=1}^K p_k^{a_k} \prod_{j=1}^q \tilde{p}_j^{b_j} = \prod_{k=1}^K p_k^{a_k} \prod_{j=1}^q (\delta_j^\top \mathbf{p})^{b_j}.$$

- $\mathbf{p} = (p_1, \dots, p_K)^\top$ : parameters of the incomplete multinomial model.
- $\mathbf{a} = (a_1, \dots, a_K)^\top$ : counts of fully classified observations.
- $\mathbf{b} = (b_1, \dots, b_q)^\top$ : counts of incomplete observations. Positive terms for partial classification, and negative terms for truncated outcomes.
- $\Delta = \{\Delta_{kj}\}_{K \times q} = [\delta_1, \dots, \delta_q]$ : indicator matrix.

# What is Incomplete Multinomial Model (Cont'd)

$$\begin{aligned} L(\mathbf{p}) &\propto p_1^{21} p_2^{24} p_3^{18} p_4^{20} p_5^{25} p_6^{12} \\ &\quad \times (p_1 + p_2 + p_3)^{18} (p_4 + p_5 + p_6)^{22} \\ &\quad \times (p_1 + p_4)^{10} (p_2 + p_5)^{20} (p_3 + p_6)^{10} \\ &\quad \times \left( \frac{p_1}{p_1 + p_4} \right)^{53} \left( \frac{p_4}{p_1 + p_4} \right)^{47}. \end{aligned}$$

$$\mathbf{a}^T = \begin{bmatrix} p_1 & p_2 & p_3 & p_4 & p_5 & p_6 \\ 21 + 53, & 24, & 18, & 20 + 47, & 25, & 12 \end{bmatrix},$$

$$\mathbf{b}^T = \begin{bmatrix} 1 & 2 & 3 & 4 & 5 \\ 18, & 22, & 10 - 53 - 47, & 20, & 10 \end{bmatrix},$$

$$\Delta^T = \begin{matrix} & j & p_1 & p_2 & p_3 & p_4 & p_5 & p_6 \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{bmatrix} 1 & 1 & 1 & & & & \\ & 1 & & 1 & 1 & 1 & \\ & & 1 & & 1 & & \\ & & & 1 & & 1 & \\ & & & & 1 & & 1 \end{bmatrix} \end{matrix}.$$

# Optimality condition

Let  $s = \sum_{k=1}^K a_k + \sum_{j=1}^q b_j$ ,  $Q^+ = \{j \mid b_j > 0, j = 1, \dots, q\}$  and  $Q^- = \{j \mid b_j < 0, j = 1, \dots, q\}$  be the sets of indices of positive and negative elements in  $\mathbf{b}$  respectively.

$$\ell(\mathbf{p}|\mathbf{a}, \mathbf{b}, \mathbf{\Delta}) = \sum_{k=1}^K a_k \log p_k + \sum_{j=1}^q b_j \log \delta_j^\top \mathbf{p} - s \left( \sum_{k=1}^K p_k - 1 \right).$$

Optimality condition:  $\nabla \ell(\mathbf{p}) = 0$ ,

$$\frac{\partial \ell}{\partial p_k} = \frac{a_k}{p_k} + \sum_{j \in Q^+} \frac{|b_j| \Delta_{kj}}{\delta_j^\top \mathbf{p}} - \sum_{j \in Q^-} \frac{|b_j| \Delta_{kj}}{\delta_j^\top \mathbf{p}} - s = 0,$$

which is equivalent to

$$a_k + \left( \sum_{j \in Q^+} \frac{|b_j| \Delta_{kj}}{\delta_j^\top \mathbf{p}} - \sum_{j \in Q^-} \frac{|b_j| \Delta_{kj}}{\delta_j^\top \mathbf{p}} - s \right) p_k = 0.$$

# Stable Weaver Algorithm

---

## Algorithm 1 Stable Weaver Algorithm

---

**Input:** Observations  $(\mathbf{a}, \mathbf{b}, \Delta)$

**Initialize:**  $\mathbf{p}^{(0)} = (1/K, \dots, 1/K)^\top$ ,  $s = \mathbf{1}^\top \mathbf{a} + \mathbf{1}^\top \mathbf{b}$

**repeat**

$\boldsymbol{\tau} = \mathbf{b} / \Delta^\top \mathbf{p}^{(t)}$  (element-wise division)

$\boldsymbol{\tau}^+ = \max(\boldsymbol{\tau}, \mathbf{0})$ ,  $\boldsymbol{\tau}^- = \min(\boldsymbol{\tau}, \mathbf{0})$

$\mathbf{p}^{(t+1)} = [\mathbf{a} + (\Delta \boldsymbol{\tau}^+) \circ \mathbf{p}^{(t)}] / (s \mathbf{1} - \Delta \boldsymbol{\tau}^-)$

( $\circ$  represents element-wise product)

$\mathbf{p}^{(t+1)} = \mathbf{p}^{(t+1)} / \text{sum}(\mathbf{p}^{(t+1)})$

**until** convergence

---

- The weaver algorithm updates the parameter by  $\mathbf{p} = \mathbf{a} / (s \mathbf{1} - \Delta \boldsymbol{\tau})$ .
- Bayesian weaver is time-consuming due to the inner-outer iteration structure and the selection of the thickening parameter is difficult.

- Contingency tables with merged and truncated cells.
- Polytomous response data with underlying categories. For example, the phenotype expressions on blood types.
- Interval censored time-to-event data with truncation in survival analysis.
- Include several well-known ranking models as special cases, such as the Bradley–Terry, Plackett–Luce models and their variants.

# Results on Real Datasets

Algorithm		NASCAR		HKJC1416	
		(w/o ties)	(w/ ties)	(w/o ties)	(w/ ties)
Stable	Iteration	22	459	40.4K	27.2K
	Weaver	Time (s)	<0.01	0.03	38.46
Bayesian	Iteration	128K	263K	>1M	>1M
	Weaver	Time (s)	25.27	50.12	>5000
MM	Iteration	22	–	40.4K	–
	Time (s)	<0.01	–	375.79	–
Trust Region*	Iteration	1937	5048	636 <sup>†</sup>	649 <sup>†</sup>
	Time (s)	74.31	125.68	1139.14	1835.37
ILSR	Iteration	12	–	4056	–
	Time (s)	0.06	–	1166.97	–
Self Consistency	Iteration	36798	11282	– <sup>‡</sup>	–
	Time (s)	11.61	2.08	–	–

\* The number of iterations for the trust region constrained algorithm refers to the number of the objective function evaluations.

<sup>†</sup> We use the approximated Hessian matrix when fitting the trust region constrained algorithm to the HKJC1416 data because its calculation is too time-consuming.

<sup>‡</sup> For the HKJC1416 data, the self-consistency approach converges to a wrong solution.



# Results on Real Datasets (Cont'd)

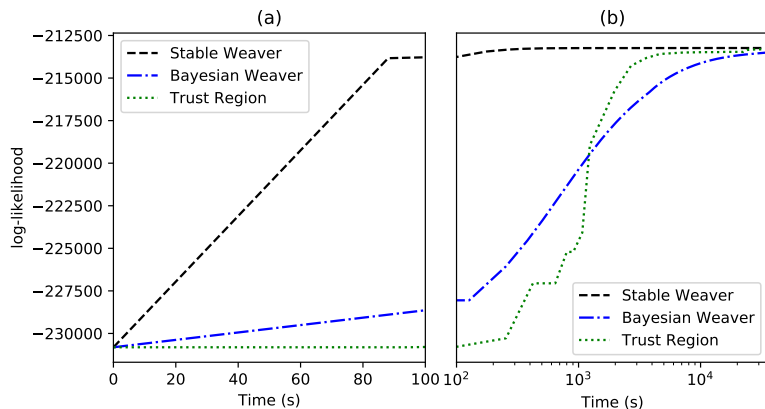


Figure 1: Convergence plot of the stable weaver algorithm compared with existing methods on the dataset HKJC9916 against running time (a)  $t \in [0, 100]$  and (b)  $t \in [100, 36000]$  (s).

Thanks for listening.