



Rank-R Multiway Logistic Regression

Fabien Girka, Pierrick Chevaillier, Arnaud Gloaguen, Giulia Gennari, Ghislaine Dehaene-Lambertz, Laurent Le Brusquet, Arthur Tenenhaus

► **To cite this version:**

Fabien Girka, Pierrick Chevaillier, Arnaud Gloaguen, Giulia Gennari, Ghislaine Dehaene-Lambertz, et al.. Rank-R Multiway Logistic Regression. 52èmes Journées de Statistique, 2021, Nice, France. hal-03051752

HAL Id: hal-03051752

<https://hal-centralesupelec.archives-ouvertes.fr/hal-03051752>

Submitted on 10 Dec 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

RANK-R MULTIWAY LOGISTIC REGRESSION

Fabien Girka ¹, Pierrick Chevaillier ¹, Arnaud Gloaguen ^{2,3}, Giulia Gennari ⁴, Ghislaine Dehaene-Lambertz ⁴, Laurent Le Brusquet ² & Arthur Tenenhaus²

¹ *CentraleSupélec, 91190, Gif-sur-Yvette, France,*

² *Université Paris-Saclay, CNRS, CentraleSupélec, Laboratoire des signaux et systèmes, 91190, Gif-sur-Yvette, France.*

³ *UNATI, Neurospin, Institut Joliot, CEA, Université Paris-Saclay, France,*

⁴ *INSERM, UMR992, Neurospin, Institut Joliot, CEA, Université Paris-Saclay, France*
fabien.girka@supelec.fr, pierrick.chevaillier@supelec.fr,
arnaud.gloaguen@centralesupelec.fr

Résumé. Beaucoup de données ont une structure intrinsèque tensorielle lorsque, par exemple, plusieurs modalités de la même variable ont été observées sur chaque individu. Les approches multivoie deviennent alors un choix naturel pour analyser ces données. Les versions standards de ces approches consistent à imposer au vecteur de poids d'être une décomposition de PARAFAC de rang 1. Pour certaines applications, les données peuvent cependant s'avérer trop complexes pour que cette décomposition soit valide. Ce papier présente une version de la régression logistique multivoie associée à une décomposition de rang R , l'objectif étant de proposer une méthode de classification applicable aux situations où la contrainte de rang 1 serait trop restrictive. Un algorithme de directions alternées est proposé pour la régression logistique multivoie de rang R . Les performances de cette méthodes sont évaluées sur des données d'électroencéphalogrammes (EEG).

Mots-clés. Analyse multivoie, régression logistique, EEG

Abstract. Data often has an inherent tensor structure (e.g. data where the same set of variables is collected at different occasions). To deal with such data, multiway models become a natural choice. Standard multiway models impose weight vectors to be rank-1 PARAFAC decomposition. However in some applications, this constraint appears to be too restrictive. This paper presents a more general version of multiway logistic regression (MLR) associated to a rank- R decomposition. The objective of such an approach is to propose a classification model that can cope with situations where rank-1 constraint may be too restrictive. An alternating direction algorithm is proposed for rank- R MLR and its performances are evaluated on electroencephalogram (EEG) data.

Keywords. Multiway analysis, logistic regression, EEG

1 Introduction

Multiway data appears in many research fields as neuroscience, chemometrics or social networks to name a few. It occurs when the same set of variables is collected through

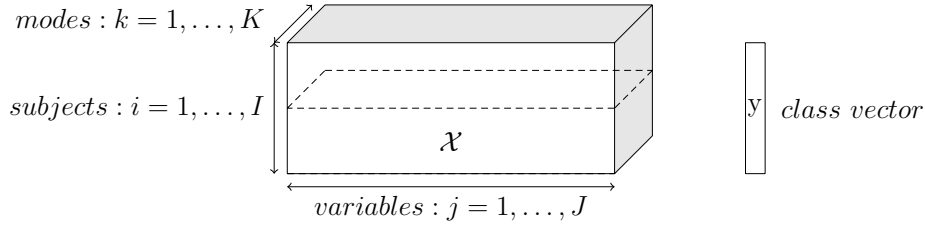


Figure 1: Three-way data, each sample is represented by K vectors of length J

different modes. For example, this is the case for spatio-temporal data (several images collected at different time steps) or when a set of measures is acquired through multiple sensors. See [1] for an overview of methods for the analysis of multiway data.

Let $\mathcal{X} = \{X_{ijk}\}_{1 \leq i \leq I, 1 \leq j \leq J, 1 \leq k \leq K}$ be a third order tensor of dimension $I \times J \times K$ where I is the number of subjects, J the number of variables and K the number of modes (see Figure 1). The roles of variables and modes are symmetrical and thus can be interchanged.

A first non-multiway option to deal with such tensor data is to unfold the tensor by concatenating its frontal slices $X_{..k}$ next to each other, leading to $\mathbf{X} = [X_{..1} \ \dots \ X_{..K}]$ which is a matrix of size $I \times JK$ on which classic statistical methods can be applied. However, such an approach leads to issues: (i) A problem of size JK could be computationally impractical for standard computers, (ii) the higher the dimension/order of the tensor, the more the number of variables obtained by flattening, the higher the risk of overfitting is, (iii) the results are not easy to interpret as the considered model does not permit a separate interpretation of the influence of the variables and modes.

A second approach largely used in the multiway literature when the model to construct uses the tensor variables through a linear form is to impose a Kronecker constraint on the weight vector: $\boldsymbol{\beta} = \boldsymbol{\beta}^J \otimes \boldsymbol{\beta}^K$. Taking into account the multiway structure of the data with a Kronecker constraint reduces the degree of freedom from JK to $J + K$, which may limit overfitting effects and computation time. Moreover, the study of the contributions of the variables and modes is made easier by analysing each vector $\boldsymbol{\beta}^J$ and $\boldsymbol{\beta}^K$ separately.

This rank-1 constraint was used for example in Multiway Logistic Regression (MLR) [2]. However, such a constraint can be too restrictive when effect of variables and modes are not strictly parallel. Therefore a higher rank decomposition for $\boldsymbol{\beta}$ can be considered :

$$\boldsymbol{\beta} = \sum_{r=1}^R \boldsymbol{\beta}_r^K \otimes \boldsymbol{\beta}_r^J \quad (1)$$

Similar rank R approaches has been studied in [3] for Generalized Linear Model and [4] for Support Vectors Machines. In this paper, we propose a rank- R Multiway Logistic Regression.

This paper is organized as follows: Section 2 presents a short reminder of (regularized) logistic regression. The rank- R logistic regression is described in Section 3. Section 4 presents an application on EEG data.

2 Regularized Logistic Regression

Logistic regression can be directly used on tensor data by unfolding the tensor in a matrix. Let \mathbf{x}_i be the vector of the KJ observed variables of subject i (\mathbf{x}_i is then the i^{th} line of the unfolded tensor \mathbf{X}), y_i its class (either 0 or 1). Logistic regression relies on maximising the conditional log-likelihood $\sum_{i=1}^n \log \mathbb{P}(y_i | \mathbf{x}_i)$ under the assumption that the conditional probabilities log-ratio is linear:

$$\log \left(\frac{\mathbb{P}(y = 1 | \mathbf{x}_i)}{1 - \mathbb{P}(y = 1 | \mathbf{x}_i)} \right) = \beta_0 + \boldsymbol{\beta}^\top \mathbf{x}_i$$

with β_0 and $\boldsymbol{\beta}$ parameters of the model. From this model it comes the following expression of the regularized log-likelihood:

$$\mathcal{C}(\beta_0, \boldsymbol{\beta}, \mathbf{X}, \mathbf{y}, \lambda) = \sum_{i=1}^n y_i (\beta_0 + \boldsymbol{\beta}^\top \mathbf{x}_i) - \log(1 + \exp(\beta_0 + \boldsymbol{\beta}^\top \mathbf{x}_i)) - \lambda g(\boldsymbol{\beta}) \quad (2)$$

where $\lambda > 0$ is a regularization parameter that can be tuned and $g(\boldsymbol{\beta})$ is a penalty term.

3 Multiway Logistic Regression

In this section, rank- R multiway logistic regression (R-MLR) is presented. R-MLR is defined as the following optimization problem:

$$\max_{\beta_0, \boldsymbol{\beta}^K, \boldsymbol{\beta}^J} \mathcal{C}(\beta_0, \boldsymbol{\beta}, \mathbf{X}, \mathbf{y}, \lambda) \quad \text{s.t.} \quad \boldsymbol{\beta} = \sum_{r=1}^R \boldsymbol{\beta}_r^K \otimes \boldsymbol{\beta}_r^J \quad (3)$$

where $\boldsymbol{\beta}^J = [(\boldsymbol{\beta}_1^J)^\top \dots (\boldsymbol{\beta}_R^J)^\top]^\top$ and $\boldsymbol{\beta}^K = [(\boldsymbol{\beta}_1^K)^\top \dots (\boldsymbol{\beta}_R^K)^\top]^\top$. An alternating direction algorithm that monotonically converges is proposed to solve the optimization problem (3). First, we can note that $\boldsymbol{\beta}^\top \mathbf{x}_i$ can be expressed as:

$$\left(\sum_{r=1}^R \boldsymbol{\beta}_r^J \otimes \boldsymbol{\beta}_r^K \right)^\top \mathbf{x}_i = \sum_{r=1}^R (\boldsymbol{\beta}_r^J)^\top ((\boldsymbol{\beta}_r^K)^\top \otimes \mathbf{I}_J) \mathbf{x}_i \doteq \sum_{r=1}^R (\boldsymbol{\beta}_r^J)^\top \mathbf{z}_{r,i}^J \quad (4)$$

In addition, two types of regularisation are considered in this paper. First an ℓ_2 penalty with $g(\boldsymbol{\beta}) = \frac{1}{2}\|\boldsymbol{\beta}\|_2^2$ that can be expressed in terms of $\boldsymbol{\beta}^J$ and $\boldsymbol{\beta}^K$ as follows:

$$\begin{aligned}\boldsymbol{\beta}^\top \boldsymbol{\beta} &= \left(\sum_{r=1}^R \boldsymbol{\beta}_r^J \otimes \boldsymbol{\beta}_r^K \right)^\top \left(\sum_{r=1}^R \boldsymbol{\beta}_r^J \otimes \boldsymbol{\beta}_r^K \right) = \sum_{r=1}^R \|\boldsymbol{\beta}_r^J\|_2^2 \|\boldsymbol{\beta}_r^K\|_2^2 + 2 \sum_{i=1}^R \sum_{j=i+1}^R (\boldsymbol{\beta}_i^J)^\top \boldsymbol{\beta}_j^J (\boldsymbol{\beta}_i^K)^\top \boldsymbol{\beta}_j^K \\ &= (\boldsymbol{\beta}^J)^\top \mathbf{R}^J \boldsymbol{\beta}^J\end{aligned}$$

where $\mathbf{R}^J = ((\mathbf{B}^K)^\top \mathbf{B}^K) \otimes \mathbf{I}_J$, with $\mathbf{B}^K = [\boldsymbol{\beta}_1^K \ \dots \ \boldsymbol{\beta}_R^K]$. As long as the columns of \mathbf{B}^K are not colinear, \mathbf{R}^J is symmetric positive definite and $\mathbf{Q}_{\ell_2}^J = (\mathbf{R}^J)^{-\frac{1}{2}} = ((\mathbf{B}^K)^\top \mathbf{B}^K)^{-\frac{1}{2}} \otimes \mathbf{I}_J$ is well defined.

Furthermore, by setting $g(\boldsymbol{\beta}) = \sum_{r=1}^R \|\boldsymbol{\beta}_r^J \otimes \boldsymbol{\beta}_r^K\|_1$, R -MLR with variable selection can also be defined. The interest of this structured sparsity-inducing norms is that, with a single parameter λ , the sparsity will spread (driven by the data) among ranks, variables and modes. This penalty term can be expressed as a function of $\boldsymbol{\beta}^J$ and $\boldsymbol{\beta}^K$:

$$\sum_{r=1}^R \|\boldsymbol{\beta}_r^J \otimes \boldsymbol{\beta}_r^K\|_1 = \|\mathbf{R}^J \boldsymbol{\beta}^J\|_1 \quad (5)$$

with $\mathbf{R}^J = (\|\boldsymbol{\beta}_r^K\|_1 \mathbf{I}_J)_{r \in \{1 \dots R\}}$ a block diagonal matrix. As before, $\mathbf{Q}_{\ell_1}^J = (\mathbf{R}^J)^{-1} = (\|\boldsymbol{\beta}_r^K\|_1^{-1} \mathbf{I}_J)_{r \in \{1 \dots R\}}$ is well defined.

As a consequence, the objective function of the optimization problem (3) can be expressed in terms $\boldsymbol{\beta}^J$ and $\boldsymbol{\beta}^K$ and therefore can be maximized with respect to β_0 , $\boldsymbol{\beta}^J$ and $\boldsymbol{\beta}^K$ using an alternating direction algorithm. Indeed, optimising w.r.t. $(\beta_0, (\mathbf{Q}^J)^{-1} \boldsymbol{\beta}^J)$ can be seen as applying logistic regression to maximise criterion $\mathcal{C}^J = \mathcal{C}(\beta_0, (\mathbf{Q}^J)^{-1} \boldsymbol{\beta}^J, \mathbf{Q}^J \mathbf{Z}^J, \mathbf{y}, \lambda)$ with \mathbf{Z}^J a matrix of size $I \times J$ and $\mathbf{z}_i^J = [(z_{1,i}^J)^\top \ \dots \ (z_{R,i}^J)^\top]^\top$. As variables and modes play symmetric roles, the same can be done with $\boldsymbol{\beta}^K$. We can now derive the R -MLR algorithm presented in Algorithm 1.

4 Application on EEG data and discussion

The objective of this study was to identify whether the infant's brain encodes the phonetic features used by linguists to describe speech. 24 different consonant-vowel syllables were presented to 25 infants in a randomized order every 1000 ms during one-hour-long experimental sessions. Brain responses were recorded at 500 Hz with a high-density EEG net comprising 252 channels. After pre-processing, this EEG experiment yields 25 tensors of size 24 syllables \times 500 time steps \times 252 channels each. The consonants varied along the manner of articulation separating the 24 syllables in 2 classes that we want to predict.

Algorithm 1: Rank R Multiway Logistic Regression

Inputs: $\epsilon > 0$, λ , R , $\beta^{K(0)}$, penalty

 $q \leftarrow 0$
repeat

$$Z_r^J = \sum_{k=1}^K \left(\beta_r^{K(q)} \right)_k X_{..k} \text{ for } r \in \{1, \dots, R\}$$

$$\left(Z^J, \beta^{K(q)} \right) \leftarrow \left(\left[(Z_1^J)^\top \dots (Z_R^J)^\top \right]^\top, \left[(\beta_1^{K(q)})^\top \dots (\beta_R^{K(q)})^\top \right]^\top \right)$$

if ℓ_1 penalty $Q^J \leftarrow Q_{\ell_1}^J$ **else** $Q^J \leftarrow Q_{\ell_2}^J$

$$(\beta_0^{(q)}, (Q^J)^{-1} \beta^{J(q)}) \leftarrow \operatorname{argmax}_{\beta_0, \beta} \mathcal{C}(\beta_0, (Q^J)^{-1} \beta, Q^J Z^J, y, \lambda)$$

$$Z_r^K = \sum_{j=1}^J \left(\beta_r^{J(q)} \right)_j X_{.j} \text{ for } r \in \{1, \dots, R\}$$

$$\left(Z^K, \beta^{J(q)} \right) \leftarrow \left(\left[(Z_1^K)^\top \dots (Z_R^K)^\top \right]^\top, \left[(\beta_1^{J(q)})^\top \dots (\beta_R^{J(q)})^\top \right]^\top \right)$$

if ℓ_1 penalty $Q^K \leftarrow Q_{\ell_1}^K$ **else** $Q^K \leftarrow Q_{\ell_2}^K$

$$(\beta_0^{(q)}, (Q^K)^{-1} \beta^{K(q+1)}) \leftarrow \operatorname{argmax}_{\beta_0, \beta} \mathcal{C}(\beta_0, (Q^K)^{-1} \beta, Q^K Z^K, y, \lambda)$$

 $q \leftarrow q + 1$
until $|\mathcal{C}^K - \mathcal{C}^J| < \epsilon \mathcal{C}^J$;

return $(\beta^{K(q)}, \beta^{J(q)}, \beta_0^{(q)})$

Regularized logistic regression and R-MLR with R from 1 to 3 are evaluated and compared using Leave-One-Out cross validation: the model is trained on all infants except one and tested on the remaining one. The Area Under the ROC Curve (AUC), the regularisation parameter λ and the computational time used to run the 25 folds are reported in Table 1.

From Table 1, we show that, for complex data such as EEG data, higher rank MLR enable to outperform both rank-1 MLR and regularized logistic regression. For ℓ_1 penalty, finding the best regularization parameter for logistic regression is really challenging given the required computation time. MLR enables to cut down drastically this computation time and yields better results for rank 3.

Model	AUC	λ	Time (in s)	Model	AUC	λ	Time (in s)
LR	0.822 ± 0.19	1000	892	LR	0.830 ± 0.19	20	117000
1-MLR	0.815 ± 0.23	6000	282	1-MLR	0.816 ± 0.22	6.5	257
2-MLR	0.852 ± 0.17	17500	340	2-MLR	0.826 ± 0.18	10	409
3-MLR	0.857 ± 0.18	17500	498	3-MLR	0.853 ± 0.19	15	557

Table 1: Cross validation results by Leave-One-Out for ℓ_2 (left) and ℓ_1 (right) penalties

In addition, multiway models enable the graphical display of the weights and get insights into the importance of variables and modes separately (see Figure 2).

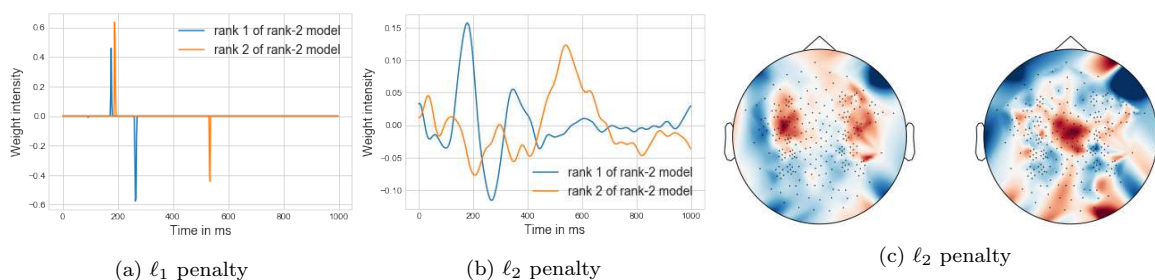


Figure 2: Weight visualisation for 2-MLR trained on all 25 infants: (a) and (b) show time step weights and (c) shows the topomap of electrode weights.

Topomaps associated with the ℓ_1 penalty (not displayed here due to lack of space) are consistent with the ℓ_2 penalty results. However, similarly to the time weights of the ℓ_1 penalty, they are very sparse with few (less than 5%) isolated selected channels. This variable selection leads spatial/time resolution that is lower than the usual phenomenon observed and described by neuroscientists. Hence, future works include to combine ℓ_1 and ℓ_2 penalties in order to try to catch smoother effects in time and space.

5 Conclusion

Rank-R Multiway Logistic Regression is presented in this paper and shows promising results in EEG application. While R-MLR is presented for third order tensors, it can be generalized to any higher order.

Bibliography

- [1] Bro, R. (2000), Multi-way Analysis in the Food Industry - Models, Algorithms, and Applications ICSLP Proceedings.
- [2] Le Brusquet L., Lechuga G., Tenenhaus A. (2014), Régression Logistique Multivoie, 46ème Journée de Statistique.
- [3] Zhou, Hua; Li, Lexin; Zhu, Hongtu (2013), Tensor Regression with Applications in Neuroimaging Data Analysis, Journal of the American Statistical Association, vol.108.
- [4] T. Lyu, E. F. Lock, and L. E. Eberly (2017), Discriminating sample groups with multi-way data, Biostatistics (Oxford, England)