

Low-rank, Orthogonally Decomposable Tensor Regression with Application to Visual Stimulus Decoding of fMRI Data

Journal:	<i>Journal of Computational and Graphical Statistics</i>
Manuscript ID	JCGS-20-152
Manuscript Type:	Original Article
Keywords:	High Dimension, Low Sample Size, Imaging Data, Low Rank Approximation, Nuclear Norm, Tensor Decomposition

SCHOLARONE™
Manuscripts

Low-rank, Orthogonally Decomposable Tensor Regression with Application to Visual Stimulus Decoding of fMRI Data

J.C. Poythress
Jeongyoun Ahn
and
Cheolwoo Park

Department of Statistics, University of Georgia

May 2, 2020

Abstract

We consider the problem of fitting a generalized linear model with a three-dimensional image covariate, such as one obtained by functional magnetic resonance imaging (fMRI). A major challenge for fitting such a model is that the image is a multidimensional array, called a tensor, containing tens of thousands of elements, called voxels. Because there is a parameter associated with each voxel, fitting the model entails estimating tens of thousands of parameters with a typical sample size on the order of hundreds of subjects. We propose to reduce the dimensionality of the problem by imposing a low rank assumption on the parameter tensor, which not only reduces the number of parameters to estimate, but also exploits the implicit spatial information in fMRI data. In addition, we assume the parameter tensor is orthogonally decomposable, enabling us to penalize the so-called tensor “singular values” to obtain a low rank estimate. We develop algorithms based on projected gradient descent and the proximal gradient method to estimate the model. In simulation, the proposed methods yielded estimates with smaller squared errors than an existing method based on alternating minimization. For visual stimulus decoding of a real fMRI dataset, the proposed methods resulted in better classification accuracies than the existing method. Visualization of the estimates revealed compact clusters of voxels with relatively large values, potentially representing brain regions relevant to the task.

Keywords: High dimension low sample size, Imaging data, Low rank approximation, Nuclear norm, Tensor decomposition.

1 1 Introduction

2
3
4
5
6 Modern technology allows researchers to collect data that can be represented as images.
7 For example, images such as those obtained by structural magnetic resonance imaging
8 (MRI), diffusion tensor imaging (DTI), electroencephalography (EEG), and functional MRI
9 (fMRI) are frequently used by psychologists, neuroscientists, and others within health and
10 medical disciplines to diagnose or study diseases. The need to analyze such imaging data
11 poses new challenges for statisticians because many traditional statistical methods do not
12 readily accommodate imaging data. Different approaches have been proposed depending
13 on the type of question the researcher would like to answer. Our work considers research
14 questions in which the image may be considered a covariate, and one would like to know
15 how the image covariate relates to or affects a response variable, which may be discrete
16 or continuous. Some examples include whether an image obtained by fMRI can be used
17 to predict whether a person would be diagnosed with a particular psychiatric disorder
18 (such as schizophrenia), or whether a structural MRI can be used to predict disease status
19 (e.g., Alzheimer's). In either example, the response may be binary (1=disorder/disease,
20 0=healthy) or continuous (e.g., symptom scale scores). Determining whether an image can
21 be used to make such a prediction and what elements of the image are important to the
22 prediction can help researchers identify the neural circuitry involved in the disorder/disease
23 and potentially contribute to a better understanding of its etiology.

24 From the statistical perspective, images can be treated as matrix-valued (two-dimensional)
25 or array-valued (three-dimensional or higher) data. However, images are challenging to
26 model using traditional methods because of the implicit spatial structure and large number
27 of parameters to estimate. For example, if one wanted to model the relationship between
28 disease status and an fMRI covariate using a generalized linear model (GLM), one would
29 need to estimate $40 \times 48 \times 38 = 72,960$ parameters – the dimension of a typical image ob-
30 tained by fMRI. One solution is to fit a regularized GLM with a sparsity-inducing penalty,
31 which effectively reduces the number of parameters to estimate. However, the results from
32 penalized regression are often unsatisfactory when the data are ultra-high dimensional.
33 Moreover, many of the classical penalized regression techniques do not account for the
34

1
2
3 spatial structure inherent in imaging data.

4
5 Several researchers have proposed tensor regression models to handle array-valued imaging
6 data. “Tensor” simply means a multidimensional array (e.g., a 1D tensor is a vector, a
7 2D tensor is a matrix, etc.). Tensor regression models exploit the expected spatial dependency
8 of the array-valued image by assuming the corresponding tensor of model parameters
9 can be well-approximated by a low rank structure. For example, if there are features that
10 distinguish? distinguish images belonging to healthy subjects from images belonging to subjects diag-
11 nosed with a disorder, we might expect the features to be few and to occur in spatially
12 compact locations. The presence of stronger signal in a few, spatially compact locations is
13 indicative of low rank structure in the overall signal.
14

15 Tensor regression models make use of tensor decompositions. Tensor decompositions
16 have a long history of applications in the fields of psychometrics, chemometrics, and sig-
17 nal processing (Kolda and Bader, 2009). In those applications, the objective is often to
18 approximate an observed tensor by one of lower rank. In contrast, for applications in
19 statistics and machine learning, the tensor to approximate by one of lower rank is usually
20 not observed. In a tensor regression model, the parameter tensor is assumed to admit
21 a low rank structure. That is, the problem is to simultaneously estimate the parameter
22 tensor and approximate it by a tensor of lower rank. Thus, the problem of estimating and
23 decomposing an unobserved quantity is of a fundamentally different nature than finding
24 an approximation to an observed quantity.
25

26 Various types of tensor regression models have been proposed. We focus on the case of
27 a scalar response and tensor covariate. Many of the models rely on the canonical polyadic
28 decomposition (CPD) (Hitchcock, 1927). The CPD decomposes a tensor as a sum of rank-1
29 tensors. For illustration, the CPD decomposes the D -dimensional tensor \mathcal{B} as
30

$$31 \quad \mathcal{B} = \sum_{r=1}^R \boldsymbol{\beta}_{1r} \circ \boldsymbol{\beta}_{2r} \circ \cdots \circ \boldsymbol{\beta}_{Dr},$$

32 where \circ denotes the vector outer product. The minimum number of terms in the sum
33 needed to reconstruct the tensor exactly is called the tensor rank. When the number of
34 terms R is truncated so that R is less than the actual tensor rank, the CPD yields a rank- R
35 approximation of the tensor. For convenience of notation, we may also collect the vectors
36

comprising the outer product into matrices such that $B_1 = (\boldsymbol{\beta}_{11}, \dots, \boldsymbol{\beta}_{1R}), \dots, B_D = (\boldsymbol{\beta}_{D1}, \dots, \boldsymbol{\beta}_{DR})$. The matrices B_1, \dots, B_D are referred to as “factor matrices.” We may then write $\mathcal{B} = [[B_1, \dots, B_D]]$.

For classification problems, Hung and Wang (2012) and Tan et al. (2013) studied logistic regression with a matrix and tensor covariate, respectively. Hung and Wang (2012) formulated the systematic part of the model as a bilinear form in the parameters, which is equivalent to a rank-1 CPD. Tan et al. (2013) considered the more general case of tensors of arbitrary dimension and a rank- R CPD. Guo et al. (2012) studied linear regression for continuous responses and support vector machine for binary responses. In both cases, they used the CPD with an additional ridge penalty or group sparsity penalty. Guhaniyogi et al. (2017) constructed a Gaussian linear regression model under a Bayesian framework by specifying priors for the vectors comprising the rank-1 terms in the CPD.

The present work focuses on tensor regression in the context of GLMs. Zhou et al. (2013) assumed a rank- R CPD for the parameter tensor in a GLM. The systematic component of their model can be expressed as

$$g(\mu_i) = \alpha + \boldsymbol{\gamma}^T \mathbf{z}_i + \langle \mathcal{B}, \mathcal{X}_i \rangle \text{ such that } \mathcal{B} = \sum_{r=1}^R \boldsymbol{\beta}_{1r} \circ \boldsymbol{\beta}_{2r} \circ \dots \circ \boldsymbol{\beta}_{Dr}, \quad (1)$$

where $g(\mu_i)$ is the link function for the mean of the response, α is a constant parameter, \mathbf{z} is a vector of regular covariates with corresponding parameter $\boldsymbol{\gamma}$, and \mathcal{X} is a tensor covariate (e.g., an fMRI image) with corresponding parameter tensor \mathcal{B} . The tensor inner product $\langle \cdot, \cdot \rangle$ is equivalent to vectorizing the two tensors and taking the vector inner product. The model assumes that the true parameter tensor can be well-approximated by a rank- R CPD. The rank- R assumption reduces the number of parameters associated with the tensor covariate from $\prod_{d=1}^D p_d$ to $R \sum_{d=1}^D p_d$, where p_d denotes the size along one dimension. For a typical fMRI image, a rank-1 assumption would reduce the number of parameters from $40 \times 48 \times 38 = 72,960$ to $40 + 48 + 38 = 126$.

The model (1), which we refer to as the CPD model, requires *a priori* specification of the rank R , so the estimate obtained is a fixed-rank approximation of the parameter. In the case of two-dimensional images, the parameter is a matrix, so the fixed-rank approximation can be thought of as hard thresholding the matrix singular values. Zhou and Li (2014)

proposed a convex relaxation of the matrix rank by penalizing the objective function with the ℓ_1 norm of the singular values, which they called spectral-regularized matrix regression. Their approach involves soft thresholding the singular values, which yields a low rank estimate of the parameter matrix by shrinking some of its singular values to exactly zero. The objective function for their model can be written as

$$\min_B -\ell\ell(B) + P_\lambda(\sigma(B)), \quad (2)$$

where B is the (unrestricted, full rank) parameter matrix, $-\ell\ell(\cdot)$ is the negative log-likelihood of a GLM with link function $g(\mu_i) = \alpha + \gamma^T z_i + \langle B, X_i \rangle$, $\sigma(\cdot)$ extracts the singular values of a matrix, and $P_\lambda(\cdot)$ is a sparsity-inducing penalty function. For example, if $P_\lambda(\cdot)$ is the LASSO penalty (Tibshirani, 1996), then $P_\lambda(\sigma(B))$ is the nuclear norm of B , $\|B\|_* = \sum_{r=1}^{\min(p_1, p_2)} \sigma_r$, multiplied by the tuning parameter λ . One of the main advantages of model (2) over model (1) is convexity. The estimates for both models must be obtained numerically, and the non-convexity of (1) makes it challenging to obtain an estimate near the global minimizer of the objective function.

The major limitation of model (2) is that it can only be fitted for two-dimensional image covariates. In Section 2, we propose two versions of a low-rank, orthogonally decomposable tensor regression model that can be applied to three-dimensional (or higher) images. Both models utilize a notion of tensor “singular values.” One version of the model assumes a fixed tensor rank, which serves as a necessary intermediate step toward developing a second version of the model that relaxes the fixed-rank assumption and instead penalizes the tensor singular values. The penalized version of the model is similar in principle to (2). However, a key difference is that matrices always admit an orthogonal decomposition (by the singular value decomposition), whereas higher-order tensors are not guaranteed to be orthogonally decomposable (“odeco” for short). We illustrate some properties of the proposed model and compare to Zhou et al.’s CPD model in a simulation experiment in Section 3. We apply the proposed methods to intra-subject visual stimulus decoding of a real fMRI dataset in Section 4. We provide some concluding remarks in Section 5.

2 Proposed Methodology

2.1 Background: An Orthogonal Tensor Decomposition

Both versions of the proposed orthogonally decomposable tensor regression model make use of a tensor decomposition studied by Chen and Saad (2009) called a low rank, orthogonal approximation to tensors (LROAT). The LROAT decomposition solves

$$\min_{\sigma_1, \dots, \sigma_R, B_1, \dots, B_D} \left\| \mathcal{B} - \sum_{r=1}^R \sigma_r \boldsymbol{\beta}_{1r} \circ \dots \circ \boldsymbol{\beta}_{Dr} \right\|_F \quad \text{s.t. } B_1^T B_1 = \dots = B_D^T B_D = I_R \quad (3)$$

$B_1 = (\beta_{11}, \dots, \beta_{1R}) \in \mathbb{R}^{d \times R}$ is an orthogonal matrix

and $\sigma_r \geq 0, r = 1, \dots, R,$

where $\|\cdot\|_F$ denotes the Frobenius norm (defined in an analogous way for tensors as it is for matrices). Chen and Saad (2009) proposed a block coordinate descent algorithm to solve (3); see their manuscript for the details.

Some important properties of the LROAT decomposition include:

1. The factor matrices B_1, \dots, B_D are orthogonal.

2. The tensor rank is bounded by the smallest size among the modes:

$$R = \text{rank}(\mathcal{B}) \leq \min(p_1, p_2, \dots, p_D).$$

3. The “singular values” are non-negative and ordered: $0 \leq \sigma_1 \leq \sigma_2 \leq \dots \leq \sigma_R.$

4. The tensor rank equals the number of nonzero singular values:

$$R = \#\{\sigma_r : \sigma_r > 0, r = 1, \dots, \min(p_1, p_2, \dots, p_D)\}.$$

why?

5. The decomposition is unique up to signs, even when some singular values are repeated.

Many of the properties listed above are shared with the matrix singular value decomposition (SVD). The key difference between LROAT for higher-order tensors and the matrix SVD is

that not all tensors are odecorable. From property 1, a tensor with $\text{rank}(\mathcal{B}) > \min(p_1, p_2, \dots, p_D)$ does not admit an orthogonal decomposition. Moreover, even if $\text{rank}(\mathcal{B}) \leq \min(p_1, p_2, \dots, p_D)$, there is no guarantee that \mathcal{B} is odecorable. Thus, in real data applications, the LROAT decomposition is nearly always an approximation of the original tensor (hence “approximation” is part of the acronym).

2.2 Fixed, Low-rank Orthogonally Decomposable Tensor Regression

To obtain the fixed-rank version of our low-rank, orthogonally decomposable tensor regression model, which we refer to as the LODTR model, we incorporate the LROAT decomposition as an additional assumption in a GLM. The systematic component of the model can be specified as

$$\begin{aligned} g(\mu_i) &= \alpha + \boldsymbol{\gamma}^T \mathbf{z}_i + \langle \mathcal{B}, \mathcal{X}_i \rangle \text{ s.t. } \text{rank}(\mathcal{B}) = R \text{ and odec} \\ &= \alpha + \boldsymbol{\gamma}^T \mathbf{z}_i + \left\langle \sum_{r=1}^R \sigma_r \boldsymbol{\beta}_{1r} \circ \boldsymbol{\beta}_{2r} \circ \cdots \circ \boldsymbol{\beta}_{Dr}, \mathcal{X}_i \right\rangle \text{ s.t. } B_1^T B_1 = \cdots = B_D^T B_D = I_R. \end{aligned} \quad (4)$$

The model (4) is similar to the model (1), except that the factor matrices are restricted to be orthogonal. In fact, when $D = 2$ (i.e., the tensors are matrices), the two models are essentially equivalent in the sense that the estimated parameters $\widehat{\mathcal{B}}$ should be equal. However, the estimated factor matrices $\widehat{B}_1, \widehat{B}_2$ will not be equal because of the orthogonality constraint. The interpretation is that the models (1) and (4) are estimating the same parameter \mathcal{B} , but the model (1) allows any basis for the column space of \mathcal{B} , while the model (4) requires an orthogonal basis for the column space of \mathcal{B} . When $D \geq 3$, the models are not equivalent because the parameter space of model (1) is all rank- R tensors of size $p_1 \times \cdots \times p_D$, while the parameter space of model (4) is all rank- R odec tensors of size $p_1 \times \cdots \times p_D$.

not all tensor are rank-R odec

We propose to estimate model (4) with a projected gradient descent algorithm. For a general problem $\min f(x)$ s.t. $x \in \mathcal{S}$, where \mathcal{S} denotes a set of constraints, the updates of the projected gradient descent algorithm take the form

$$x^{(k+1)} = \prod_{\mathcal{S}} (x^k - \delta^k \nabla f(x^k)), \quad \text{how to calculate?}$$

where δ^k is a step size and $\prod_{\mathcal{S}}$ denotes the projection operator onto the feasible set \mathcal{S} . The projection operator can be defined as $\prod_{\mathcal{S}}(x) = \arg \min_{z \in \mathcal{S}} \|x - z\|_F$. The projected gradient descent method can be interpreted as the usual gradient descent method with the additional modification that each update is required to lie in the feasible set. The projection step ensures that each update lies in the feasible set while minimizing the distance between the update and the update that would be obtained from an unconstrained version of the problem.

The constraint set in model (4) is the space of rank- R odec tensors. To project onto that set, recall problem (3). Problem (3) is identical to the problem defined by projecting onto the space of rank- R odec tensors. This suggests that the algorithm described by Chen and Saad (2009) may be used for the projection step in a projected gradient descent algorithm for fitting model (4). The gradient descent step may be based on a suitable loss function, such as the negative log-likelihood. The details of the algorithm for fitting model (4) are included in the Supplementary Material.

2.3 Penalized Orthogonally Decomposable Tensor Regression

To modify model (4) to avoid *a priori* specification of the rank, we penalize the objective function implied by the model with a function of the tensor singular values. The penalized version of our model, which we refer to as penalized orthogonally decomposable tensor regression (PODTR), can be written as

$$\min_{\mathcal{B}} -\ell\ell(\mathcal{B}) + P_\lambda(\sigma(\mathcal{B})), \quad (5)$$

where \mathcal{B} is the parameter tensor (assumed to be odec), $-\ell\ell(\cdot)$ is the negative log-likelihood of a GLM, $\sigma(\cdot)$ extracts the singular values of an odec tensor, and $P_\lambda(\cdot)$ is a sparsity-inducing penalty function. For this study, we only consider the LASSO penalty [i.e., $P_\lambda(\sigma(\mathcal{B}))$ is the ℓ_1 norm of the tensor singular values, multiplied by λ]. By Property 4 of the LROAT decomposition (Section 2.1), the LASSO penalty applied to the tensor singular values will yield a low rank estimate of \mathcal{B} .

$\sum \sigma_i < a \Rightarrow \sigma_n, \dots, \sigma_m > 0,$
 $\sigma_{m-1}, \dots, \sigma_1 = 0$, where m depends on λ .

To estimate model (5), we adapt Zhou and Li's algorithm for solving (2), which is based on the proximal gradient method. The standard proximal gradient method applied to (2) splits the problem into two steps: one that minimizes the negative log-likelihood by gradient descent, and one that applies the proximal operator associated with $P_\lambda(\sigma(\mathcal{B}))$. The proximal operator associated with the LASSO penalty is the soft thresholding operator, so solving (2) involves shrinking the singular values of B to zero. Zhou and Li's algorithm modifies the standard proximal gradient method by incorporating an extrapolation step due to Nesterov (1983) that achieves an accelerated convergence rate.

We must further modify the proximal gradient method to apply it to (5). Namely, we must account for the fact that tensors of order $D \geq 3$ must be projected onto the set of odecor tensors before shrinking the singular values. We propose to include a step that projects onto the set of rank-min(p_1, \dots, p_D) odecor tensors – the largest rank possible for the LROAT decomposition of a tensor of size $p_1 \times \dots \times p_D$. When the actual tensor rank $R < \min(p_1, \dots, p_D)$, this does not pose a problem because we will have $\sigma_{R+1} = \dots = \sigma_{\min(p_1, \dots, p_D)} = 0$. A brief sketch of the updates for the algorithm is

$$\text{Gradient descent step: } \dot{\mathcal{B}}^{(k+1)} = \mathcal{B}^{(k)} - \delta^k \nabla f(\mathcal{B}^{(k)})$$

$$\begin{aligned} \text{Projection step: } \ddot{\mathcal{B}}^{(k+1)} &= \arg \min \left\| \dot{\mathcal{B}}^{(k+1)} - \sum_{r=1}^{\min(p_1, \dots, p_D)} \sigma_r \beta_{1r} \circ \dots \circ \beta_{Dr} \right\|_F \\ &\text{s.t. } B_1^T B_1 = \dots = B_D^T B_D = I_{\min(p_1, \dots, p_D)} \end{aligned}$$

$$\text{Soft thresholding step: } \mathcal{B}^{(k+1)} = \sum_{r=1}^{\min(p_1, \dots, p_D)} (\ddot{\sigma}_r - \delta^k \lambda)_+ \ddot{\beta}_{1r} \circ \dots \circ \ddot{\beta}_{Dr},$$

after update, we may get a low-structure with different rank. How to ensure the rank will converge?

where the function f denotes the negative log-likelihood and $(\cdot)_+ = \max(0, \cdot)$. The complete details of the proposed algorithm for fitting model (5) are included in the Supplementary Material.

my own question: a rank 6 cp structure tensor can be best approximated by a rank 5 cp structure? Does this fact have some influence on the inference of rank?

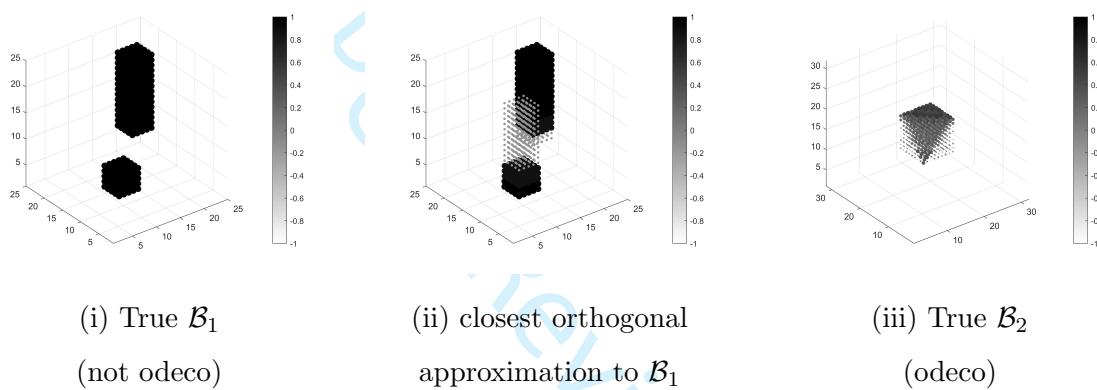
3 Simulation Experiment

We used simulated data to illustrate some of the properties of the LODTR and PODTR models compared to the CPD model. For $i = 1, \dots, N$, we generated the response as $Y_i \sim \text{Normal}(\beta_0 + \langle \mathcal{X}_i, \mathcal{B} \rangle, 1)$ and the “image” covariate $(\mathcal{X}_i)_{jkl} \sim \text{Normal}(0, 1)$. We set $\beta_0 = 1$ in all of our experiments.

We compare what happens when the true parameter tensor \mathcal{B} is odecor vs. not. We constructed several different parameter tensors:

1. A $25 \times 25 \times 25$ rank-2 non-odeco tensor in which the signal region was comprised of two hyperrectangles [labelled \mathcal{B}_1 in Fig. 1(i)]. The values of the signal elements were set to 1 and all other elements were set to 0.

1
2
3 2. A $32 \times 32 \times 32$ rank-4 non-odeco tensor with two signal regions: one comprised of a
4 cross and the other comprised of nested cubes (see the Supplementary Material for the
5 visualization). The values of the signal elements of the cross, inner cube and outer cube
6 were set to 1, 1.5, and 0.75, respectively, with all other elements set to 0.
7
8 why there is no comparison between: rank 4 non-odeco vs rank 4 odec
9
10 3. A $32 \times 32 \times 32$ rank-9 odec tensor using the singular values and corresponding singular
11 vectors from the SVD of a two-dimensional image of a triangle [labelled \mathcal{B}_2 in Fig. 1(iii)].
12 See the accompanying MATLAB code in the Supplementary Material for additional
13 details about constructing the odec tensor.
14
15
16
17



34 Figure 1: True parameter tensors for the simulation experiment.
35
36

We simulated data from the models 100 times and calculated the squared error of the estimates, defined as $\|\mathcal{B} - \hat{\mathcal{B}}\|_F^2$. The quartiles of the squared errors are shown in Table 1. The LODTR and PODTR estimates typically have smaller squared errors than the CPD estimates, even when the true parameter tensor is not odec. Most of the poor performance of the CPD model can be attributed to its sensitivity to the initial value, as evidenced by the large discrepancies between Q1 and Q3. For the hyperrectangles parameter, the first quartile of the squared errors of the rank-2 CPD estimates was much lower than those for LODTR and PODTR, indicating that the relaxed assumptions of the CPD model can be advantageous, but only if the initialization happens to be good (which seems to be infrequently).

59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32	31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1
----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	---	---	---	---	---	---	---	---	---

Table 1: Summary of the squared errors of the estimated parameter tensors. Quartiles are calculated based on the squared errors from 100 simulations; the best value in each row is highlighted in bold. The sample size N for each simulation is given in the first column.

		CPD sensitive to the initial points					where is the rank result of PODTR?											
		CPD					LODTR					PODTR						
		Rank					Rank					Tuning Parameter λ						
		1	2	3	4	5	1	2	3	4	5	2500	2000	1500	1000	250	100	
\text{\\tX_i \\in R^{32,}\\ 32,32}, i \\in [N].	Hyperrectangles (rank-2, non-odeco) $N = 400$	Q1	152.9	0.6	404.5	1323.7	2480	152.3	18.1	30.7	45.3	56.6	123.4	93.6	67.2	54.1	45.1	46.1
		Q2	157.8	254	1344.4	2090.6	3017.9	155.9	18.9	33.1	48	60.3	129.8	98.2	73.4	59.4	50.7	52
		Q3	168	964.5	1660.4	2371.3	3329.4	159.6	19.6	35.4	52	68	139.2	107.1	80.3	64.4	55.1	58
11	Cross + nested cubes (rank-4, non-odeco) $N = 500$	Q1	76.6	131.9	525.1	859.5	1070.1	73.1	31.5	54.4	79.8	96.5	77.4	67.7	62.5	59.6	60.2	64.9
		Q2	277.2	419.9	639.5	929.7	1200.8	75.1	32.7	58.7	86.7	102	82	72	66.6	63.5	67.2	74.3
		Q3	291.4	447.9	691.5	1032.6	1364.4	77.2	34.4	62.9	92	111.4	86	77.3	71.4	69.2	76.8	98.8
11	Triangle (rank-9, odec) $N = 500$	Q1	11.2	20.1	38.3	229.4	290.8	11	9.5	17.4	24.2	29	31.5	25	19.2	15.1	14.6	15.5
		Q2	11.8	21.9	169.5	257.3	316.9	11.3	10.9	19.4	28	33.3	33	26.1	19.8	16.1	16	16.6
		Q3	77.8	122.5	182.5	275.6	352.7	11.6	18	27.1	34.3	39.5	35.2	27.9	21.2	17.1	16.8	17.6
11	Triangle (rank-9, odec) $N = 1000$	Q1	9.8	6.3	7.9	10.5	15.3	9.8	6	4.5	4.3	5.6	15.7	13	10	7.3	6.3	7.4
		Q2	9.9	12.3	8.6	12.1	25.1	9.9	6.1	4.6	5	6.7	16.3	13.4	10.3	7.6	6.8	8.2
		Q3	10.1	12.8	16.5	21.3	27.4	10	6.3	4.7	6.3	8.2	16.7	13.9	10.7	8	7.3	9.4

the best performances in higher rank correspond to a smaller rank than true rank

If we visualize the estimates, several additional properties of the LODTR and PODTR models vs. the CPD model become apparent. Figure 1 shows the true hyperrectangle and triangle parameters (labelled \mathcal{B}_1 and \mathcal{B}_2 , respectively), while Figures 2 and 3 show several instances of their corresponding estimates (see the Supplementary Material for additional estimates). For illustration, we also show the closest orthogonal approximation to the hyperrectangle parameter [Fig. 1(ii)], which we obtained by applying Chen and Saad's algorithm. Notice that the closest orthogonal approximation captures the main features of the true signal, but there are some unavoidable artifacts caused by the orthogonality constraint. Comparing the true parameters to their estimates, we can observe that:

1. The CPD model performs well when the sample size is large relative to the number of parameters to estimate, but breaks down quickly as the number of parameters approaches the sample size [e.g., Fig. 2(iii) and Fig. 3(iii)].
2. When the true parameter tensor is not odec, the LODTR and PODTR estimates improve as higher rank models are fitted, but reach a point beyond which no further improvement is possible [e.g., Fig. 2(iv)–(vi)]. However, the higher rank estimates closely resemble the closest orthogonal approximation of the true parameter tensor, so the estimates still manage to capture most of the main features.
3. When the true parameter tensor is odec, the LODTR and PODTR models perform better the CPD model, especially for higher rank models [e.g., Fig. 3(ii) vs. Fig. 3(v)].
4. The PODTR model is able to produce estimates that appear intermediate in rank. For example, Figure 2(vii) is intermediate between the zero tensor and the rank-1 estimates from the other models, and Figure 2(viii) is intermediate between the rank-1 and rank-2 estimates from the other models. Such a property is a result of continuous shrinkage of the singular values as opposed to fitting discrete rank models.

it is because
your true rank is
2???

why intermediate? why final estimate does not have a specific rank?
so singular value is very small?

does this mean the selection of lambda is not so good?

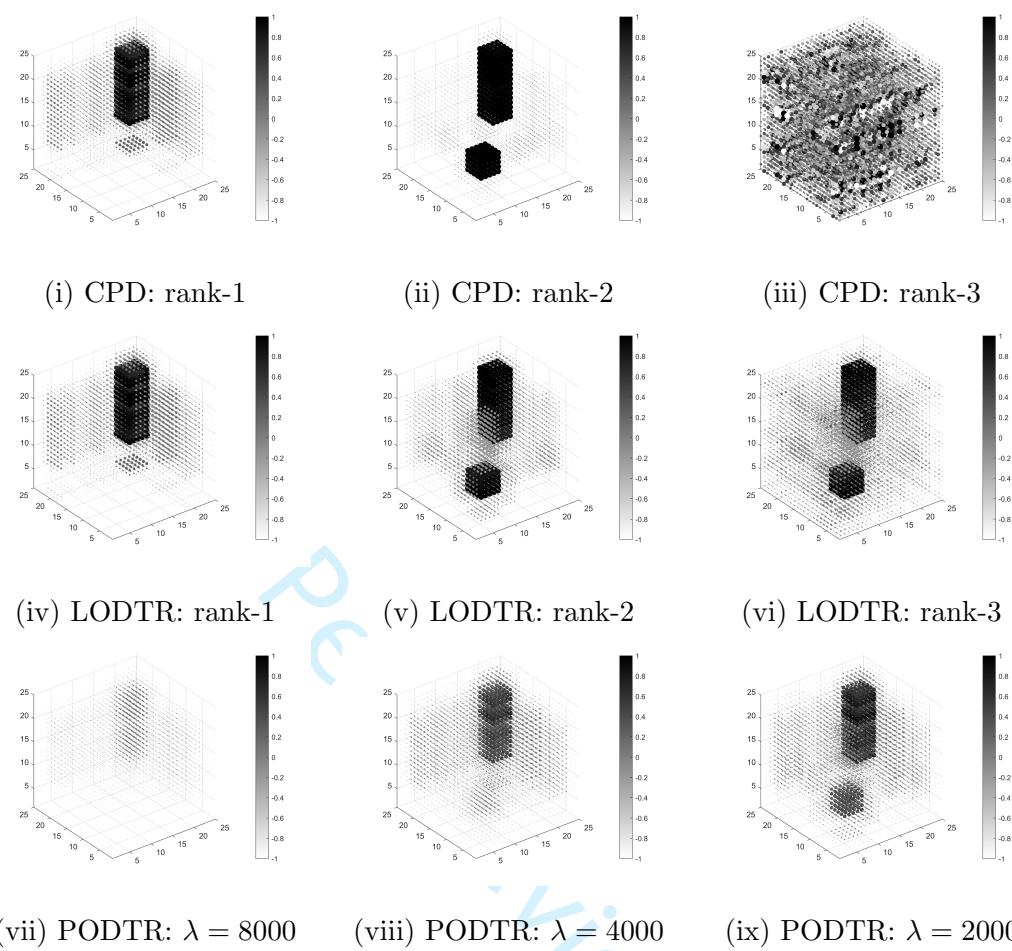


Figure 2: Estimates of the non-odeco tensor \mathcal{B}_1 under the CPD, LODTR, and PODTR models. Sample size: $N = 400$.

4 Application to Intra-subject Visual Stimulus Decoding

One application of the proposed methods is intra-subject visual stimulus decoding. One first obtains a set of fMRI images as a subject performs several visual stimulus tasks (typically, viewing different types or categories of images). Then, in the encoding step of visual stimulus decoding, one uses the fMRI images to train a classification model in which the fMRI image itself serves as the predictor and the type of image being viewed (i.e., the task) serves as the response. For the decoding step, one uses the model to predict the type of image being viewed for a separate test set of fMRI images. The goal is to identify the brain regions involved for one task versus another.

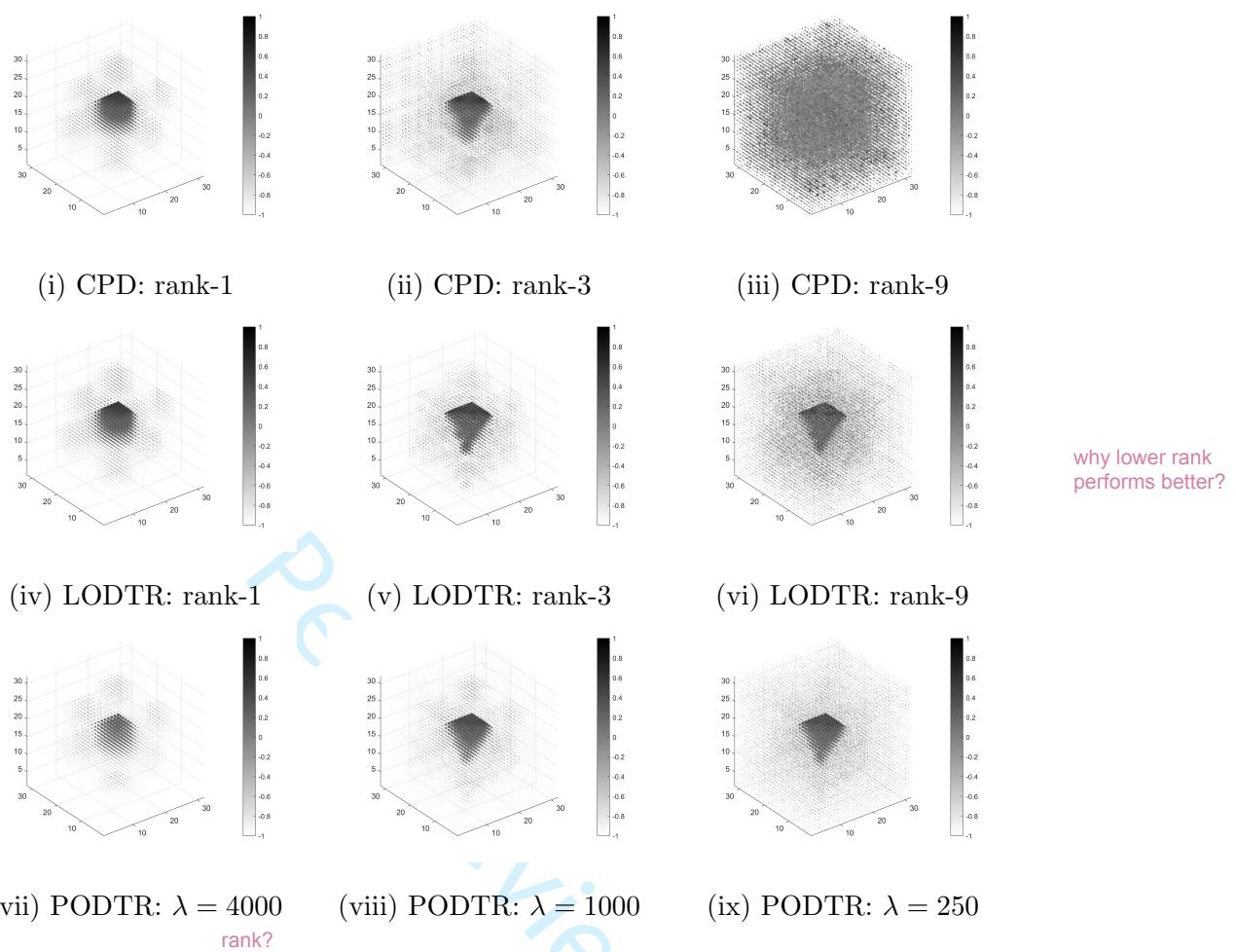


Figure 3: Estimates of the odec tensor \mathcal{B}_2 under the CPD, LODTR, and PODTR models.

Sample size: $N = 1200$.

Haxby et al. (2001) showed that the Visual Object Recognition (VOR) dataset can be used for visual stimulus decoding. The VOR data (Haxby et al., 2001; Hanson et al., 2004; O'Toole et al., 2005) is freely available through the OpenNeuro project (openneuro.org) and contains fMRI data for six subjects as they viewed different types of images (viewing each type of image is considered a task). There were eight different types of images: bottles, scissors, shoes, chairs, houses, cats, faces, and scrambled images. Data were recorded for each subject for 12 runs (except Subject 5, who had data for 11 runs). For each run, the image types were shown sequentially in blocks of 24 s duration with 12 s of rest between blocks of different image type. Within a block, images of the same type were shown every

2
3 2 s (images of the same type includes different objects of the same type and different views
4 of the same object).
5

6 The data for one run consists of two parts: 1) an events file giving the time of onset
7 of an image and the type of image and 2) a $40 \times 64 \times 64 \times 121$ fMRI tensor containing
8 raw BOLD signal. The $40 \times 64 \times 64$ component of the fMRI tensor represents the spatial
9 component with voxel dimensions $3.50\text{ mm} \times 3.75\text{ mm} \times 3.75\text{ mm}$. The component of
10 length 121 represents the time component with dimension 2.5 s . For each block of image
11 type, we extracted the time slices corresponding to 6 s after the onset of the first image
12 (to account for the delay in the hemodynamic response) up to the onset of the last image
13 in the block. Because the resolution of the time component of the fMRI tensor does not
14 match up exactly with the duration of a block, the number of time slices we obtained for
15 a block varied slightly across image types (± 1 time slice).
16

17 We applied the CPD, LODTR, and PODTR models to the VOR data for binary visual
18 stimulus decoding. We used a logistic regression model with the indicator of image type as
19 the response and the subset of the fMRI tensor corresponding to the visual cortex as the
20 predictor. The visual cortex is located in the posterior part of the brain. By visualizing a
21 whole-brain fMRI image, we determined that indices 1–40 along the first spatial dimension,
22 indices 9–30 along the second spatial dimension, and indices 23–48 along the third spatial
23 dimension should conservatively capture the visual cortex. The size of the resulting image
24 was $40 \times 22 \times 26$. We trained the models on the data for 11 runs ($N^{\text{train}} \approx 145$), then obtained
25 the classification accuracies based on the predictions for the remaining run ($N^{\text{test}} \approx 13$).
26 We repeated the process for all 6 subjects and each of their 11 or 12 runs, for a total of
27 71 fitted models per method. For clarity of exposition, we refer to the results for run j as
28 those for the models trained with the data from run j left out.
29

30 For our classification tasks, we focused on classifying scrambled images vs. each of the
31 other image types. We compared scrambled images against each of the others because the
32 signal should be strongest for those comparisons, and so they serve as good examples to
33 illustrate the uses of the proposed models. Scrambled images should result in general acti-
34 vation of the visual cortex, whereas images of other types may result in stronger activation
35
36
37
38
39
40
41
42
43
44
45

of specific subregions within the visual cortex; visualizing the parameter estimates from the trained models can help identify the subregions of the visual cortex that are associated with a particular task. For LODTR and CPD, we were restricted to fitting rank-1 models. The limited sample size relative to the dimension of the fMRI tensor precluded fitting higher rank models. For PODTR, we fitted the model for a sequence of values of the tuning parameter. For each type of image comparison, we reported the classification accuracy achieved for the test dataset. We classified an image as 1 if $\hat{p} > 0.5$ and 0 if $\hat{p} < 0.5$, where \hat{p} was defined as

$$\hat{p} = \frac{1}{1 + \exp[-\hat{\beta}_0 - \langle \mathcal{X}_i, \hat{\mathcal{B}} \rangle]}.$$

The average classification accuracies achieved by each method are shown in Table 2, where the averages are taken across both subjects and runs. Similar results disaggregated by subject can be found in the Supplementary Material. Using the fMRI image as input, all three models were able to predict which type of visual stimulus was being viewed better than by random guessing. The LODTR and PODTR models performed uniformly better than the CPD model, and sometimes substantially better. The PODTR model performed better than the LODTR model except for some instances in which a large amount of regularization was used. In general, the average classification accuracy of the PODTR model was better with less regularization, and its performance declined as more regularization was added.

In addition to making accurate predictions, another goal of the proposed models is to understand what components of the image are most related to the response. For visual stimulus decoding, that can be interpreted as understanding which areas of the brain are responsible for discriminating between one type of viewed image and another, as measured by the BOLD signal. If we visualize the estimated parameter tensor, we hope to find compact clusters of voxels with relatively large magnitude values, which we can interpret as regions of the brain that are associated with performing one task versus another.

Figures 4 – 6 show the absolute values of the average estimate of the parameter tensor over Subject 1's 12 runs for selected tasks. Results for the other tasks can be found in the Supplementary Material. The PODTR model estimates are shown for two values of the tuning parameter, representing different extremes with respect to the amount of shrinkage

Table 2: Average classification accuracy (%) for the CPD, LODTR, and PODTR models applied to the VOR data. Averages are taken across 6 subjects \times 11 or 12 runs per subject = 71 fitted models per method.

Task	CPD (rank-1)	LODTR (rank-1)	PODTR				
			$\lambda = 50$	$\lambda = 100$	$\lambda = 500$	$\lambda = 1000$	$\lambda = 2000$
scrambled vs. face	73.29	85.52	90.94	91.19	89.58	86.24	80.91
scrambled vs. house	70.62	82.84	89.41	89.96	88.51	86.72	84.25
scrambled vs. bottle	65.74	72.56	80.07	80.83	79.47	76.37	73.39
scrambled vs. scissors	62.61	71.75	83.86	83.37	82.38	81.51	75.04
scrambled vs. cat	67.05	83.14	86.03	85.48	85.49	83.41	78.18
scrambled vs. chair	64.21	77.12	80.15	81.33	80.20	78.73	73.22
scrambled vs. shoe	60.78	73.14	87.22	85.08	83.86	81.88	76.67

($\lambda = 1000$ results in large shrinkage and lower rank estimates, while $\lambda = 50$ results in less shrinkage and higher rank estimates). Note that the colorbars for the PODTR estimates are typically on a smaller scale than those for the rank-1 CPD and LODTR estimates, reflecting the powerful effect of shrinkage. The size of the dots in the PODTR estimates have been magnified to a comparable level as the CPD and LODTR estimates to enable easier identification of regions with (relatively) large values.

In all of the figures, regions of the estimated tensor with large magnitude values (relative to other regions within the same estimate) can be interpreted as signal regions – i.e., regions that are strongly associated with the type of image being viewed. The CPD estimates [Figs. 4(i)–6(i)] are noisy, with the largest values occurring in scattered voxels rather than in spatially compact locations. Such estimates are difficult to interpret because no clearly defined regions stand out. Although the CPD estimates are less interpretable, the classification accuracies were better than random guessing for all of the tasks.

The rank-1 LODTR estimates [Figs. 4(ii)–6(ii)] are less noisy than the CPD estimates, especially for the scrambled vs. face [Fig. 4(ii)] and scrambled vs. house [Fig. 5(ii)]. For those estimates, distinct, spatially compact regions of the estimate stand out, suggesting that the corresponding region of the brain is relevant to the task. For scrambled vs. bottle [Fig. 6(ii)], the voxels with relatively large magnitude values are few and they do not form

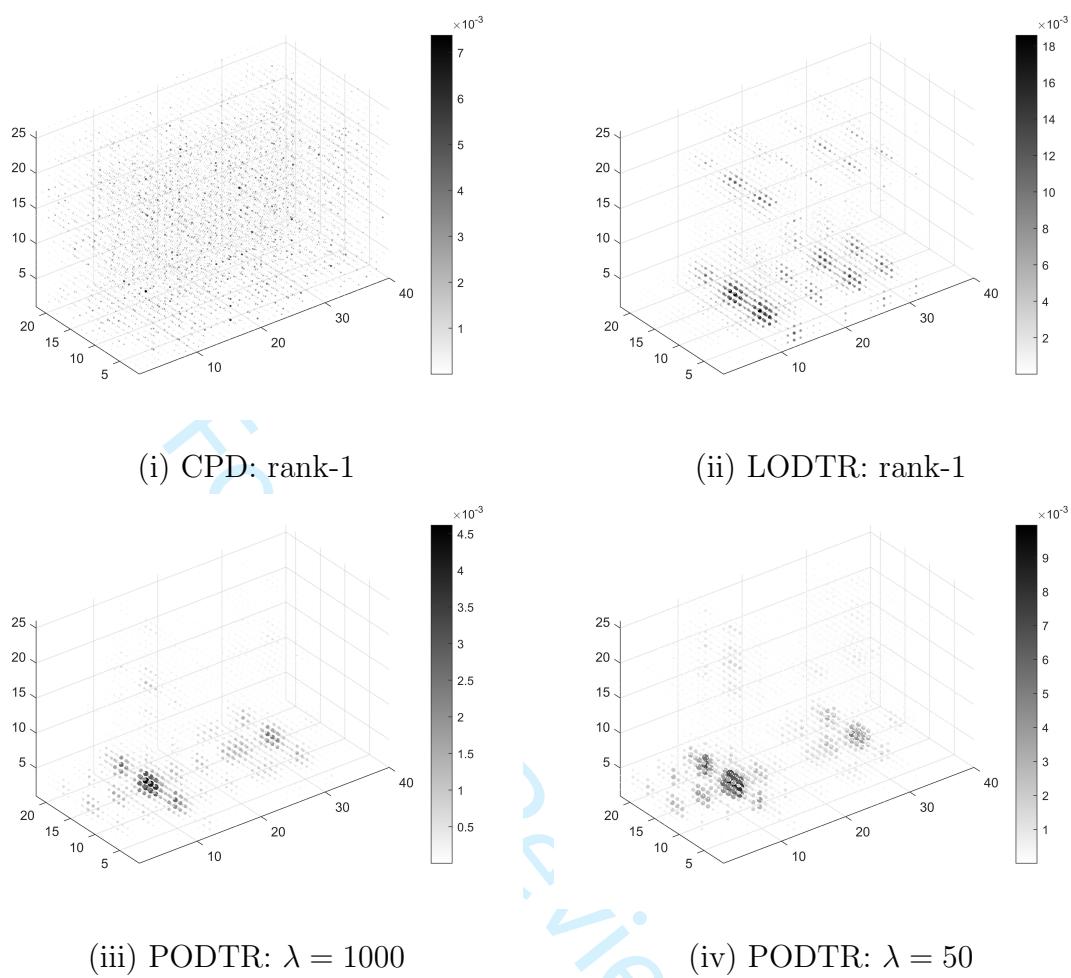


Figure 4: Scrambled vs. Face: Average $\widehat{\mathcal{B}}$ for Subject 1's 12 runs.

compact clusters, making them difficult to interpret

The PODTR estimates [Figs. 4(iii)–6(iii) and 4(iv)–6(iv)] are easiest to interpret and become more interpretable with less shrinkage. The same regions that stood out in the rank-1 LODTR estimates for scrambled vs. face and scrambled vs. house stand out in the PODTR estimates as well. The most improvement occurs for the scrambled vs. bottle task under less shrinkage. Although certain regions of the estimates stand out to some extent with greater shrinkage [Fig. 6(iii)], they are much more well-defined with less shrinkage [Fig. 6(iv)]. From Table 2, the scrambled vs. bottle estimates obtained with less shrinkage had better classification accuracies than those obtained with greater shrinkage. Thus, we can say that the estimates with large magnitude values located in compact clusters tended

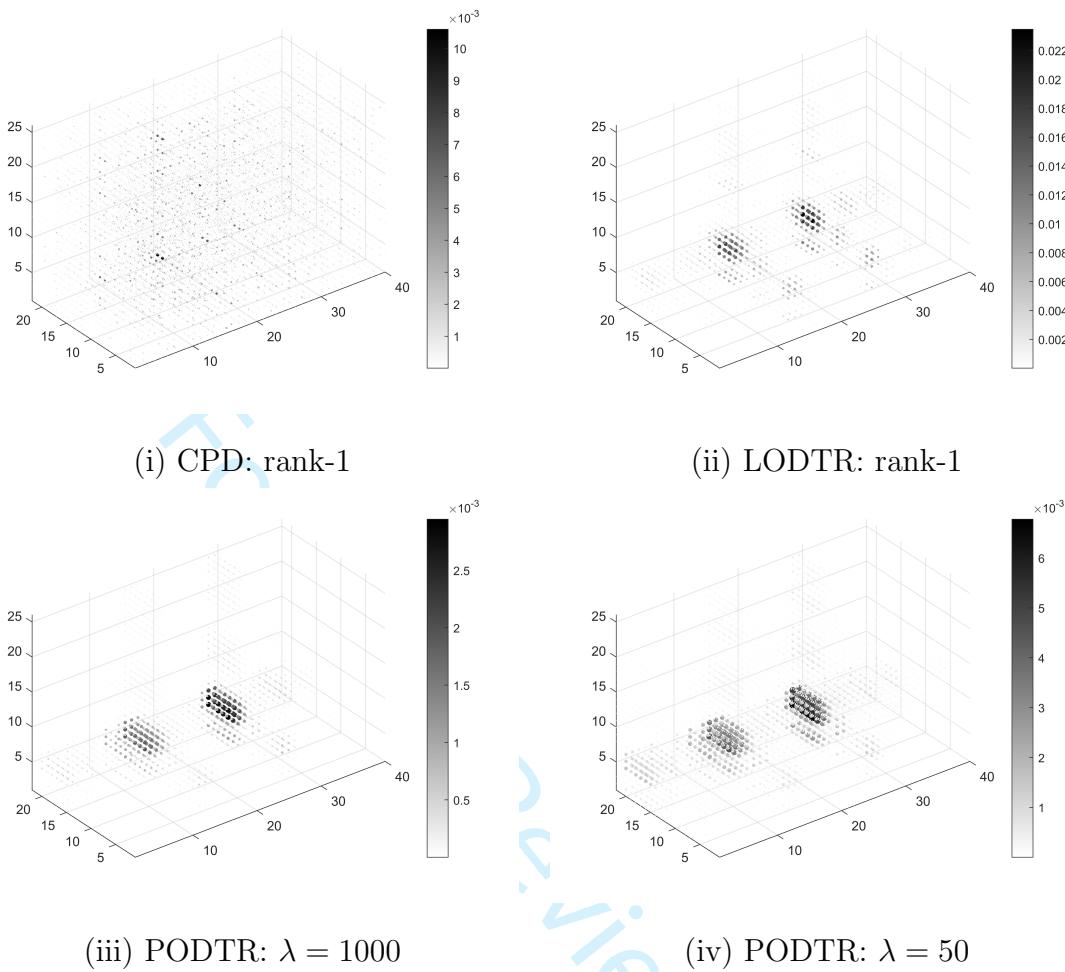


Figure 5: Scrambled vs. House: Average $\widehat{\mathcal{B}}$ for Subject 1's 12 runs.

to yield better classification accuracies for predictions made on new data.

The clusters of voxels with relatively large values in Figures 4–6 represent regions within the visual cortex associated with each of the tasks. A natural question is: Do the clusters of voxels have biological relevance? Namely, do they correspond to regions of the brain that are already known to be associated with each of the tasks? Based on the relative positions of the clusters across tasks, the clusters of voxels in Figure 4 appear to correspond to the fusiform face area, the clusters in Figure 5 to the parahippocampal place area, and the clusters in Figure 6 to the lateral occipital complex. The fusiform face area, parahippocampal place area, and lateral occipital complex have been shown to be associated with viewing images of faces, locations and things (Kanwisher et al., 2001), respectively, which is consistent

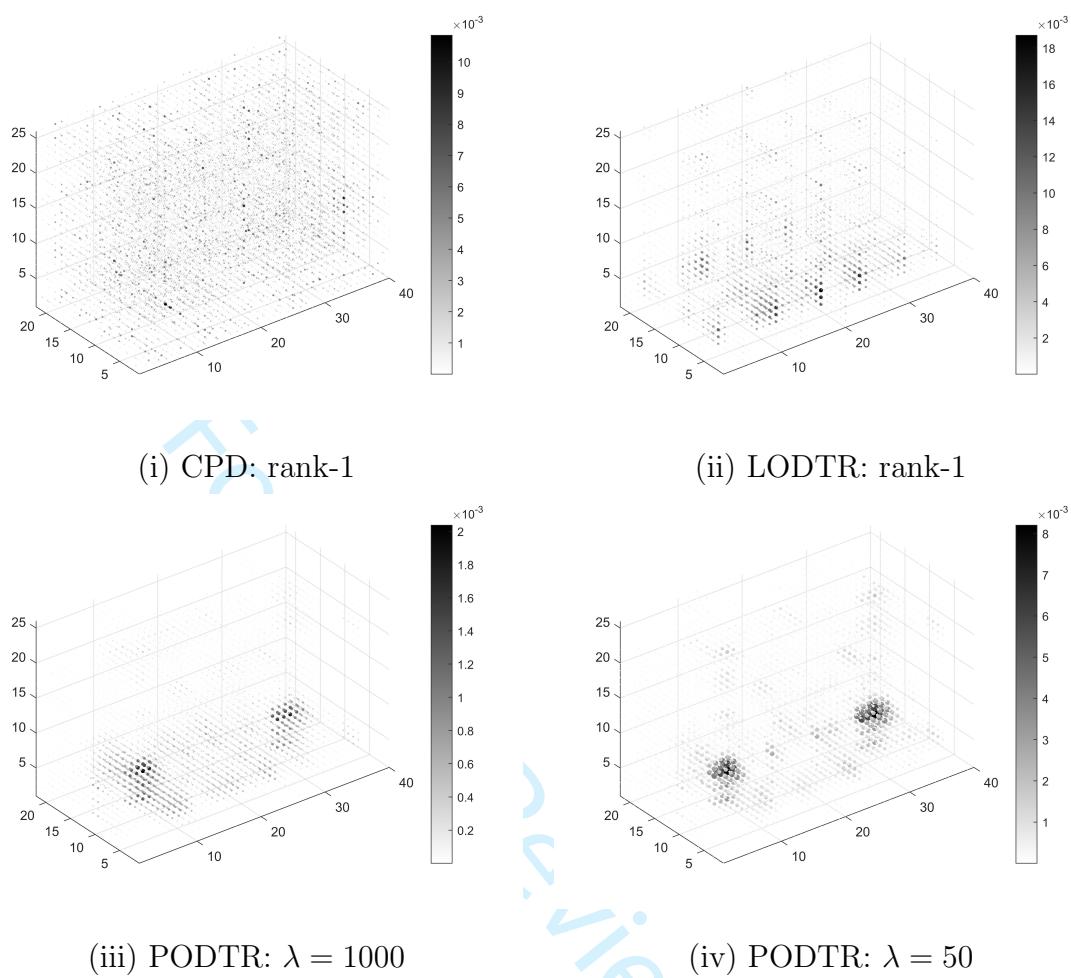


Figure 6: Scrambled vs. Bottle: Average $\widehat{\mathcal{B}}$ for Subject 1's 12 runs.

with the tasks shown. We caution that such an assessment is tentative at best. Expert judgment would be necessary to confirm that the regions corresponding to the clusters of voxels are indeed those that have been identified by others.

5 Conclusions and Future Directions

Through a simulation experiment and the analysis of real fMRI dataset, we demonstrated that the proposed low-rank, orthogonally decomposable tensor regression models have the potential to perform better than existing methods both with respect to predictive performance and interpretation. Because the LODTR model is a special case of the CPD model

to which we compared (and the two are in fact equivalent in some instances), the superior performance of the LODTR model can be attributed to some extent to the estimation algorithm rather than the model itself. The projected gradient descent algorithm used to fit the LODTR model yielded parameter estimates with smaller squared errors in simulation and higher classification accuracies in our real data analysis. In both the simulation and real data analysis, the parameter estimates tended to be easier to interpret when visualized. In contrast, the alternating minimization algorithm used to fit the CPD model yielded parameter estimates that were easy to interpret only for lower rank models in the simulation; the algorithm completely failed to produce interpretable estimates in the real data analysis and higher rank models in the simulation. Both algorithms used the same initial value, tolerance criterion, and tolerance value, so the differences cannot be attributed to differences in the initialization or stopping rule. Furthermore, the LODTR and CPD models fitted for the real data analysis were all rank 1, and in that case the LODTR and CPD models are equivalent. One possible explanation for the poor performance of the alternating minimization algorithm is that alternating among parameters makes the algorithm more likely to find local optima, especially when the number of parameters is large relative to the sample size. The projected gradient descent algorithm updates all parameters at once, perhaps making it less likely to get stuck in a local optimum.

The results from the real data application were also consistent with some of the properties of the models demonstrated in the simulation experiment. In the simulation, the performance of the CPD model was observed to break down quickly as the number of parameters approached the sample size. In the real data application, fitting the logistic regression model entailed estimating 88 parameters with only 145 observations, so it is not surprising that the CPD model produced poor estimates. The simulation also revealed that the LODTR and PODTR models are not capable of estimating the true parameter tensor perfectly when it is not odec. Real data do not follow a model, so there is no concept of a “true” parameter tensor, and even if there was, we wouldn’t know what it would look like. Then, we do not know whether the clusters of voxels evident in Figures 4–6 are actually voxels associated with the tasks, or whether some of the clusters are potentially artifacts

1
2
3 caused by the orthogonality constraints. Besides not knowing what a “good” estimate of
4 the parameter should look like, we also cannot compare with the CPD estimates (which are
5 unconstrained) because the CPD model performed so poorly. At the least, we can say that
6 the estimates with the most clearly defined clusters also produced the best classification
7 accuracies in Table 2, suggesting that the clusters of voxels are in fact associated with the
8 tasks.
9

10 In future work we hope to explore the convergence properties of the algorithms we have
11 proposed for fitting the LODTR and PODTR models. Empirically, the algorithms appear
12 to converge based on the simulated and real examples we have tried. Moreover, the solutions
13 are frequently better than the solutions from the alternating minimization algorithm
14 of Zhou et al. (2013). Unfortunately, there are myriad challenges for proving global conver-
15 gence for tensors of dimension ≥ 3 . One challenge is that both of the algorithms we have
16 proposed rely on Chen and Saad’s algorithm for the LROAT decomposition. They were
17 unable to prove global convergence of their algorithm. In addition, their algorithm is not
18 guaranteed to converge to the global optimum of problem (3). Thus, the projection steps
19 in the proposed algorithms are not necessarily carried out exactly. That creates a difficulty
20 for applying any known results about the convergence of projected gradient descent with
21 non-convex constraints. Another challenge is in dealing with the non-convex constraints
22 imposed by the low rank and orthogonality assumptions. Some results regarding the con-
23 vergence of projected gradient descent are known in the context of low rank matrices [see
24 Jain and Kar (2017)] and have been applied in the context of the Tucker decomposition
25 for tensors [e.g., Chen et al. (2019)]. However, given that not all tensors admit an or-
26 thogonal decomposition, it is not clear how one might extend those results to the LROAT
27 decomposition utilized in the proposed algorithms. Finally, a major challenge in proving
28 convergence of the algorithm for fitting for the PODTR model is that it is neither a pro-
29 jected gradient descent algorithm nor a proximal gradient algorithm in the strict sense,
30 but rather a mixture of the two. Although superficially similar to Zhou and Li’s algorithm
31 for spectral-regularized matrix regression, the proposed algorithm involves an additional
32 projection step onto the set of odec tensors. For matrices, this is simply the SVD and is a
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 change of basis rather than a projection. For tensors, it is a non-convex projection for all
4 but the rare cases in which the tensor already admits an orthogonal decomposition. Since
5 the proposed algorithm adds a non-trivial modification to the standard proximal gradient
6 method, it may not be possible to apply any known convergence results for the proximal
7 gradient method.
8
9
10
11
12
13

14 Acknowledgement

15

16
17 The third author's work was partly supported by National Science Foundation (NSF IIS-
18 1607919).
19
20
21

22 SUPPLEMENTARY MATERIAL

23

24
25 **Algorithms:** A detailed description of the algorithms used to estimate the LODTR and
26 PODTR models. (.pdf file)
27
28

29 **Additional results:** Complete set of figures and results associated with the simulation
30 experiment and VOR data analysis. We also include some computational results
31 related to the algorithms, such as the number of iterations until convergence and
32 computation time. (.pdf file)
33
34
35

36 **MATLAB code:** Collection of scripts and functions necessary to reproduce the simula-
37 tion results and real data analysis. (.m files)
38
39
40
41

42 References

43

44
45 Chen, H., G. Raskutti, and M. Yuan (2019). Non-convex projected gradient descent for
46 generalized low-rank tensor regression. *Journal of Machine Learning Research* 20(5),
47
48 1–37.

49
50
51 Chen, J. and Y. Saad (2009). On the tensor SVD and the optimal low rank orthogonal
52 approximation of tensors. *SIAM Journal on Matrix Analysis and Applications* 30(4),
53
54 1709–1734.
55
56
57
58
59
60

- 1
2
3 Guhaniyogi, R., S. Qamar, and D. B. Dunson (2017). Bayesian tensor regression. *Journal*
4
5 *of Machine Learning Research* 18(79), 1–31.
6
7 Guo, W., I. Kotsia, and I. Patras (2012, Feb.). Tensor learning for regression. *IEEE*
8
9 *Transactions on Image Processing* 21(2), 816–827.
10
11 Hanson, S. J., T. Matsuka, and J. V. Haxby (2004). Combinatorial codes in ventral tem-
12 poral lobe for object recognition: Haxby (2001) revisited: is there a “face” area? *Neu-*
13 *roImage* 23(1), 156 – 166.
14
15 Haxby, J. V., M. I. Gobbini, M. L. Furey, A. Ishai, J. L. Schouten, and P. Pietrini (2001).
16
17 Distributed and overlapping representations of faces and objects in ventral temporal
18 cortex. *Science* 293(5539), 2425–2430.
19
20
21 Hitchcock, F. L. (1927). The expression of a tensor or a polyadic as a sum of products.
22
23 *Journal of Mathematics and Physics* 6(1-4), 164–189.
24
25 Hung, H. and C.-C. Wang (2012, 07). Matrix variate logistic regression model with appli-
26 cation to EEG data. *Biostatistics* 14(1), 189–202.
27
28 Jain, P. and P. Kar (2017). Non-convex optimization for machine learning. *Foundations*
29
30 *and Trends in Machine Learning* 10(3-4), 142–363.
31
32 Kanwisher, N., P. Downing, R. Epstein, and Z. Kourtzi (2001). Functional neuroimaging of
33 visual recognition. In R. Cabeza and A. Kingstone (Eds.), *Handbook of Functional Neu-*
34 *roimaging of Cognition* (1 ed.), Chapter 5, pp. 109–152. Cambridge MA: Massachusetts
35 Institute of Technology Press.
36
37 Kolda, T. G. and B. W. Bader (2009). Tensor decompositions and applications. *SIAM*
38
39 *Review* 51(3), 455–500.
40
41 Nesterov, Y. (1983). A method of solving a convex programming problem with convergence
42 rate $O(1/k^2)$. *Soviet Mathematics Doklady* 27(2), 372–376.
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

- O'Toole, A. J., F. Jiang, H. Abdi, and J. V. Haxby (2005). Partially distributed representations of objects and faces in ventral temporal cortex. *Journal of Cognitive Neuroscience* 17(4), 580–590.
- Tan, X., Y. Zhang, S. Tang, J. Shao, F. Wu, and Y. Zhuang (2013). Logistic tensor regression for classification. In J. Yang, F. Fang, and C. Sun (Eds.), *Intelligent Science and Intelligent Data Engineering*, Berlin, Heidelberg, pp. 573–581. Springer Berlin Heidelberg.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 58(1), 267–288.
- Zhou, H. and L. Li (2014). Regularized matrix regression. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 76(2), 463–483.
- Zhou, H., L. Li, and H. Zhu (2013). Tensor regression with applications in neuroimaging data analysis. *Journal of the American Statistical Association* 108(502), 540–552.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20 Supplementary Material for “Low-rank, Orthogonally
21
22 Decomposable Tensor Regression with Application to
23
24 Visual Stimulus Decoding of fMRI Data”
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

Algorithms for Estimating the LODTR and PODTR Models

Algorithm 1: Projected Gradient Descent for LODTR Model

```

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Input: Response  $Y_{N \times 1}$ , regular covariates  $Z_{N \times p_0}$ , tensor covariate  $\mathcal{X}_{p_1 \times \dots \times p_D \times N}$ ,  

initial guesses for parameters  $\alpha^{(0)}, \gamma^{(0)}, \mathcal{B}^{(0)}$ , assumed rank  $R$ , and  

distribution (normal or Bernoulli, though any distribution in the  

exponential family in theory).  

Output: Final estimates  $\hat{\alpha}, \hat{\gamma}, \hat{\mathcal{B}} = [[\text{diag}(\hat{\sigma}_1, \dots, \hat{\sigma}_R); \hat{B}_1, \dots, \hat{B}_D]]$ .  

// Force initial guess to be rank-R and orthogonally-decomposable.  

 $[B_1 \dots, B_D] \leftarrow \text{hosvd}(\mathcal{B}^{(0)}, R_1 = \dots = R_D = R)$   

// Calls higher-order SVD algorithm described in Kolda and Bader (2009)  

 $\mathcal{B}^{(0)} \leftarrow \text{lroat}(\mathcal{B}^{(0)}, R, B_1 \dots, B_D)$   

// Calls LROAT algorithm described by Chen and Saad (2009)  

 $k \leftarrow 0$  // iteration counter  

 $\delta^0 \leftarrow 1$  // initialize step size  

while objective criterion not met do  

    // gradient descent step  

 $\text{vec}(\alpha^{(\text{temp})}, \gamma^{(\text{temp})}, \tilde{\mathcal{B}}^{(\text{temp})}) \leftarrow \text{vec}(\alpha^{(k)}, \gamma^{(k)}, \mathcal{B}^{(k)}) - \delta^k \nabla f(\alpha^{(k)}, \gamma^{(k)}, \mathcal{B}^{(k)})$   

// function  $f$  is  $-\ell\ell$   

    // projection step  

 $[B_1 \dots, B_D] \leftarrow \text{hosvd}(\tilde{\mathcal{B}}^{(\text{temp})}, R_1 = \dots = R_D = R)$   

 $\mathcal{B}^{(\text{temp})} \leftarrow \text{lroat}(\tilde{\mathcal{B}}^{(\text{temp})}, R, B_1 \dots, B_D)$   

    if  $f(\alpha^{(\text{temp})}, \gamma^{(\text{temp})}, \mathcal{B}^{(\text{temp})}) < f(\alpha^{(k)}, \gamma^{(k)}, \mathcal{B}^{(k)})$  then  

         $[\alpha^{(k+1)}, \gamma^{(k+1)}, \mathcal{B}^{(k+1)}] \leftarrow [\alpha^{(\text{temp})}, \gamma^{(\text{temp})}, \mathcal{B}^{(\text{temp})}]$  // accept update  

         $\delta^{k+1} \leftarrow 1.2 * \delta^k$  // increase step size  

    else  

        // reject update; perform line search for step size  

         $[\alpha^{(k+1)}, \gamma^{(k+1)}, \mathcal{B}^{(k+1)}] \leftarrow [\alpha^{(k)}, \gamma^{(k)}, \mathcal{B}^{(k)}]$   

         $\delta^{k+1} \leftarrow 0.5 * \delta^k$  // shrink step size  

    end  

     $k \leftarrow k + 1$  // update iteration counter  

end  

 $[\hat{\alpha}, \hat{\gamma}, \hat{\mathcal{B}}] \leftarrow [\alpha^{(\text{final})}, \gamma^{(\text{final})}, \mathcal{B}^{(\text{final})}]$ 
```

Algorithm 2: Proximal Gradient Method for PODTR Model

```

1
2
3
4
5   Input: Response  $Y_{N \times 1}$ , regular covariates  $Z_{N \times p_0}$ , tensor covariate  $\mathcal{X}_{p_1 \times \dots \times p_D \times N}$ , tuning parameter  $\lambda$ , and
6   distribution (normal, Bernoulli, or another distribution in the exponential family).
7   Output: Final estimates  $\hat{\alpha}, \hat{\gamma}, \hat{\mathcal{B}} = [[diag(\hat{\sigma}_1(\lambda), \dots, \hat{\sigma}_{\min(p_1, \dots, p_D)}(\lambda)); \hat{B}_1, \dots, \hat{B}_D]]$ .
8
9   // Note:  $\theta := vec(\alpha, \gamma, \mathcal{B})$ 
10  // initial guesses all zero
11   $\alpha^{(0)} = \alpha^{(1)} \leftarrow 0; \gamma^{(0)} = \gamma^{(1)} \leftarrow 0_{p_0}; \mathcal{B}^{(0)} = \mathcal{B}^{(1)} \leftarrow 0_{p_1 \times \dots \times p_D}$ 
12   $k \leftarrow 0$  // iteration counter
13   $\delta^0 \leftarrow 1$  // initialize step size
14   $\eta^0 \leftarrow 0; \eta^1 \leftarrow 1$  // initialize extrapolation parameters
15
16  while objective criterion not met do
17    // extrapolation step
18     $\theta^{(k)} = \theta^{(k)} + \frac{\eta^{k-1} - 1}{\eta^k} (\theta^{(k)} - \theta^{(k-1)})$ 
19
20    while not descended do
21      // gradient descent step
22       $vec(\alpha^{(temp)}, \gamma^{(temp)}, \tilde{\mathcal{A}}) \leftarrow \theta^{(k)} - \delta^k \nabla f(\theta^{(k)})$ 
23      // function  $f$  is  $-\ell$ 
24      //  $\tilde{\mathcal{A}}$  is part of an intermediate update for  $\mathcal{B}^{(k)}$ 
25
26      // projection step
27       $[A_1 \dots, A_D] \leftarrow \text{hosvd}(\tilde{\mathcal{A}}, R_1 = \dots = R_D = \min(p_1, \dots, p_D))$ 
28      // Calls higher-order SVD algorithm described in Kolda and Bader (2009)
29       $\mathcal{A} \leftarrow \text{lroat}(\tilde{\mathcal{A}}, R = \min(p_1, \dots, p_D), A_1 \dots, A_D)$ 
30      // Calls LROAT algorithm described by Chen and Saad (2009)
31
32      // thresholding step
33      for  $r = 1, \dots, \min(p_1, \dots, p_D)$  do
34         $|\sigma_r^{(temp)} \leftarrow \max(0, a_r - \delta^k \lambda)|$  //  $a_r$ 's are singular values of  $\mathcal{A}$ 
35
36      // update  $\mathcal{B}^{(k)}$  with thresholded singular values and orthogonal factor matrices of  $\mathcal{A}$ 
37       $\mathcal{B}^{(temp)} \leftarrow [[diag(\sigma_1^{(temp)}, \dots, \sigma_{\min(p_1, \dots, p_D)}^{(temp)}); A_1 \dots, A_D]]$ 
38
39      // check if descent
40      if  $f(\theta^{(temp)}) < f(\theta^{(k)}) + \nabla f(\theta^{(k)})^T (\theta^{(temp)} - \theta^{(k)}) + \frac{1}{2\delta^k} \|\theta^{(temp)} - \theta^{(k)}\|_2^2$  then
41         $[\alpha^{(k+1)}, \gamma^{(k+1)}, \mathcal{B}^{(k+1)}] \leftarrow [\alpha^{(temp)}, \gamma^{(temp)}, \mathcal{B}^{(temp)}]$  // accept update
42      else
43        // reject update; perform line search for step size
44         $\delta^k \leftarrow 0.5 * \delta^k$  // shrink step size
45      end
46
47       $\eta^{k+1} \leftarrow 0.5 * (1 + \sqrt{1 + (2\eta^k)^2})$  // update extrapolation parameter
48       $k \leftarrow k + 1$  // update iteration counter
49    end
50
51   $[\hat{\alpha}, \hat{\gamma}, \hat{\mathcal{B}}] \leftarrow [\alpha^{(final)}, \gamma^{(final)}, \mathcal{B}^{(final)}]$ 
52
53
54
55
56
57
58
59
60

```

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Bibliography

- Chen, J. and Y. Saad (2009). On the tensor SVD and the optimal low rank orthogonal approximation of tensors. *SIAM Journal on Matrix Analysis and Applications* 30(4), 1709–1734.
- Kolda, T. G. and B. W. Bader (2009). Tensor decompositions and applications. *SIAM Review* 51(3), 455–500.