Hu Sun¹ Ward Manchester² Meng Jin³ Yang Liu⁴ Yang Chen¹

Abstract

Multi-channel imaging data is a prevalent data format in scientific fields such as astronomy and biology. The structured information and the high dimensionality of these 3-D tensor data makes the analysis an intriguing but challenging topic for statisticians and practitioners. The low-rank scalar-on-tensor regression model, in particular, has received widespread attention and has been reformulated as a tensor Gaussian Process (Tensor-GP) model with multi-linear kernel in Yu et al. (2018). In this paper, we extend the Tensor-GP model by introducing an integrative dimensionality reduction technique, called tensor contraction, with a Tensor-GP for a scalar-on-tensor regression task with multi-channel imaging data. This is motivated by the solar flare forecasting problem with high dimensional multi-channel imaging data. We first estimate a latent, reduced-size tensor for each data tensor and then apply a multi-linear Tensor-GP on the latent tensor data for prediction. We introduce an anisotropic total-variation regularization when conducting the tensor contraction to obtain a sparse and smooth latent tensor. We then propose an alternating proximal gradient descent algorithm for estimation. We validate our approach via extensive simulation studies and applying it to the solar flare forecasting problem.

1. Introduction

Regression models that deal with scalar labels and tensor covariates, i.e. scalar-on-tensor regression, have received widespread attention over the past decade (Hung & Wang, 2013; Zhou et al., 2013; Zhou & Li, 2014; Kang et al., 2018;

Proceedings of the 40th International Conference on Machine Learning, Honolulu, Hawaii, USA. PMLR 202, 2023. Copyright 2023 by the author(s).

Li et al., 2018; Papadogeorgou et al., 2021). Given m-mode tensor covariate $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_m}$ and scalar label $y \in \mathbb{R}$, the existing literature approaches the regression problem mainly via:

$$\mathbb{E}[y|\mathcal{X}] = \alpha + \langle \mathcal{W}, \mathcal{X} \rangle, \tag{1}$$

where α is the intercept, \mathcal{W} is the regression coefficient tensor that matches the shape of \mathcal{X} and $\langle \cdot, \cdot \rangle$ denotes tensor inner product following Kolda & Bader (2009). This formulation can be readily adopted under the framework of generalized linear model (Zhou et al., 2013) while simultaneously preserving the tensor structure of \mathcal{X} . Typically, tensor data is of ultra-high dimensions and thus \mathcal{W} is also of high dimensionality. Various constraints have been introduced on \mathcal{W} , such as tensor norm regularization (Guo et al., 2011; Zhou & Li, 2014) and tensor rank constraints (Papadogeorgou et al., 2021; Hao et al., 2021). These constraints induce a sparse and low-rank structure over \mathcal{W} , making inferences of the high-order correlation between the scalar label and the tensor covariates tractable and interpretable.

Gaussian Process (GP) (Williams & Rasmussen, 2006) is an alternative approach to modeling complex correlation structures, and has been applied to tensor regression problems (Kang et al., 2018), where a GP prior is imposed on \mathcal{W} . In Yu et al. (2018), it is established that the tensor regression model (1), together with a low-rank constraint on \mathcal{W} , leads to the same estimator $\widehat{\mathcal{W}}$ as the tensor Gaussian Process (**Tensor-GP**) coupled with a multi-linear kernel on the prior of \mathcal{W} . A multi-linear kernel function $k(\cdot,\cdot)$ for m-mode tensors $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_m}$ can be defined in a Kronecker product form as:

$$k(\mathcal{X}_{i}, \mathcal{X}_{j}) = \mathbf{vec}\left(\mathcal{X}_{i}\right)^{\top} \left(\bigotimes_{m'=1}^{m} \mathbf{K}_{m+1-m'}\right) \mathbf{vec}\left(\mathcal{X}_{j}\right),$$

where $\mathbf{vec}\left(\cdot\right)$ is the vectorization operator and \otimes denotes the matrix Kronecker product and $\mathbf{K}_{1},\ldots,\mathbf{K}_{m}$ capture the mode-specific covariance structure of the regression coefficient tensor \mathcal{W} and are assumed to be low-rank. Interpreting this GP regression model can be hard since one needs to inspect the multi-linear kernel which deals with the tensor data at its original dimensionality $d=\prod_{m'=1}^{m}I_{m'}$.

The capability of the multi-linear Tensor-GP to provide uncertainty quantification on the prediction makes it an attractive alternative to its counterpart in (1), but a sufficient

¹Department of Statistics, University of Michigan, Ann Arbor ²Department of Climate and Space Sciences and Engineering, University of Michigan, Ann Arbor ³Solar & Astrophysics Lab, Lockheed Martin ⁴W.W. Hansen Experimental Physics Laboratory, Stanford University. Correspondence to: Yang Chen <ychenang@umich.edu>.

dimension reduction on the tensor data is needed to make it more interpretable for scientific applications. In a different thread of literature, in Kossaifi et al. (2020), a tensor contraction operation is introduced before estimating the tensor regression model under the neural network settings. Instead of compressing the information of tensor data into a vector, the tensor data is contracted into a smaller *core* tensor with the same number of modes. Such a dimension reduction technique preserves the tensor structure of the data, making tensor regression or Tensor-GP directly applicable.

In this paper, we propose a novel framework by combining the merits of tensor contraction and Tensor-GP for the scalar-on-tensor regression task. Our framework consists of two major blocks. Firstly, we introduce tensor contraction to transform the tensor data $\mathcal X$ to a feature tensor $\mathcal Z$ with much lower dimensionality. Secondly, we apply the multi-linear Tensor-GP to the reduced-sized tensor $\mathcal Z$ for regression. We build our model around a special type of tensor, i.e. the multi-channel imaging tensor, motivated by an application to astrophysical imaging analysis. But our model can be easily extended to a general tensor setup. We summarize our contributions as follows:

- We integrate tensor dimension reduction with Tensor-GP in a unified framework called Tensor-GPST, allowing for learning a low-dimensional tensor representation in a supervised learning context.
- We propose to use the anisotropic total variation regularization (Wang et al., 2017) in the tensor contraction step for a sparse and spatially smooth tensor dimension reduction. We estimate the parameters of Tensor-GPST jointly under a penalized marginal likelihood approach coupled with the proximal gradient method (Parikh et al., 2014) with convergence guarantee.

2. Tensor Gaussian Process with Spatial Transformation (Tensor-GPST)

In this section, we will first introduce our method, called Tensor-GPST, for the scalar-on-tensor regression task and then discuss the algorithm in section 2.2 for estimating its parameters and conclude by discussing the theoretical guarantee of the algorithm convergence in section 2.3.

Throughout the paper, we use calligraphic letters (e.g. \mathcal{X}, \mathcal{Z}) for tensors with at least three modes, boldface uppercase letters (e.g. \mathbf{A}, \mathbf{B}) for matrices, boldface lowercase letters (e.g. \mathbf{w}, \mathbf{y}) for vectors and plain letters (e.g. λ, s) for scalars. For an m-mode tensor \mathcal{X} of size $I_1 \times I_2 \times \cdots \times I_m$, its k-mode product with matrix $\mathbf{U} \in \mathbb{R}^{J \times I_k}$, denoted as $\mathcal{X} \times_k \mathbf{U}$, is an m-mode tensor of size $I_1 \times \ldots I_{k-1} \times J \times I_{k+1} \times I$

 $\cdots \times I_m$, where:

$$(\mathcal{X} \times_k \mathbf{U})_{i_1,\dots,j,\dots,i_m} = \sum_{i_k=1}^{I_k} \mathcal{X}_{i_1,\dots,i_k,\dots,i_m} \mathbf{U}_{ji_k}.$$

We use $\langle \mathcal{X}, \mathcal{Y} \rangle$ to denote tensor inner product and $\|\mathcal{X}\|_F = \sqrt{\langle \mathcal{X}, \mathcal{X} \rangle}$ to denote tensor Frobenius norm. We refer the readers to Kolda & Bader (2009) for a thorough introduction to tensor algebra.

2.1. Method

We consider a multi-channel imaging dataset $\{\mathcal{X}_i, y_i\}_{i=1}^N$, where $\mathcal{X}_i \in \mathbb{R}^{H \times W \times C}$ with H, W, C as the height, width and number of channels, respectively; and $y_i \in \mathbb{R}$. We use $\mathbf{X}_i^{(c)} \in \mathbb{R}^{H \times W}, c \in [C]$ to denote the c^{th} channel of \mathcal{X}_i . Gaussian process regression (GPR) (Williams & Rasmussen, 2006) specifies the prior for y_i as:

$$y_i = f(\mathcal{X}_i) + \epsilon_i, \quad f(\cdot) \sim GP(m(\cdot), k(\cdot, \cdot)), \quad (2)$$

with $\epsilon_i \sim \mathcal{N}\left(0,\sigma^2\right)$ being the idiosyncratic noise. The GP prior characterizes the unknown function $f(\cdot)$ evaluated at all data points as a multivariate Gaussian distribution, with a mean function $m(\cdot)$ and a covariance kernel function $k(\cdot,\cdot)$. Typically, $m(\cdot)$ is assumed to be zero and $k(\cdot,\cdot)$ fully specifies the behavior of the GP prior.

Given the high dimensionality of \mathcal{X}_i , it would be difficult to directly estimate and interpret the tensor kernel $k(\cdot,\cdot)$. Here we consider adding one extra step called *tensor contraction*, which compresses the information of $\mathcal{X}_i \in \mathbb{R}^{H \times W \times C}$ into a reduced-sized tensor $\mathcal{Z}_i \in \mathbb{R}^{h \times w \times C}$, with h < H, w < W, via:

$$\mathcal{Z}_i = g(\mathcal{X}_i) = \mathcal{X}_i \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{I}_C, \tag{3}$$

with $\mathbf{A} \in \mathbb{R}^{h \times H}$, $\mathbf{B} \in \mathbb{R}^{w \times W}$. In effect, \mathbf{A} and \mathbf{B} reduce the dimension of each channel of \mathcal{X}_i from $H \times W$ to $h \times w$ and one can rewrite (3) equivalently as:

$$\mathbf{Z}_i^{(c)} = \mathbf{A} \mathbf{X}_i^{(c)} \mathbf{B}^\top, \quad c = 1, 2, \dots, C.$$

After (3), we then apply (2) on \mathcal{Z}_i , as discussed later.

This formulation of tensor contraction can be found in a more general setting in tensor regression networks (Kossaifi et al., 2020), where tensor contraction can be applied to compress any tensors in a neural network. In our method, we envelope the tensor contraction operation within a tensor GP framework. Also, note that in (3), all channels share the same tensor contracting factors $\bf A$ and $\bf B$, which preserves the spatial consistency of different channels of the reduced-sized tensor $\bf Z$ for easier interpretation. Alternatively, one can replace the $\bf I_C$ in (3) with an arbitrary $\bf C \times \bf C$ matrix $\bf C$, ending up with the full tensor contraction in Kossaifi et al. (2020). We stick to (3) for simplicity in this paper.

One can interpret the contracted tensor \mathcal{Z}_i as the latent low-dimensional representation of the original tensor \mathcal{X}_i . Each $(s,t)^{\text{th}}$ element of $\mathbf{Z}_i^{(c)}$ is constructed via a matrix inner product with a rank-1 "feature map": $\mathbf{Z}_i^{(c)}(s,t) = \left\langle \boldsymbol{\alpha}_s^{\top} \boldsymbol{\beta}_t, \mathbf{X}_i^{(c)} \right\rangle$, where $\boldsymbol{\alpha}_s$ and $\boldsymbol{\beta}_t$, the basis of the feature map, are the s^{th} and t^{th} rows of \mathbf{A} and \mathbf{B} , respectively. We denote the feature map $\left(\boldsymbol{\alpha}_s^{\top} \boldsymbol{\beta}_t\right)$ as $\mathbf{W}_{s,t} \in \mathbb{R}^{H \times W}$. A visual explanation of the tensor contraction operation is shown in Figure 1a. Note how elements of $\mathbf{Z}_i^{(c)}$ on the same row or column share the same feature map basis in \mathbf{A} or \mathbf{B} .

Given the transformed tensor $\mathcal{Z}_i = g(\mathcal{X}_i)$, we assume a GP prior for $\mathbf{y} = (y_1, y_2, \dots, y_N)^{\top}$ given $\mathcal{Z}_1, \mathcal{Z}_2, \dots, \mathcal{Z}_N$ with a multi-linear kernel (Yu et al., 2018):

$$y_i = h(\mathcal{Z}_i) + \epsilon_i, \quad h(\cdot) \sim GP(0, k(\cdot, \cdot)),$$
 (4)

where $k(\cdot, \cdot)$ is the multi-linear tensor kernel function:

$$k(\mathcal{Z}_i, \mathcal{Z}_j) = \mathbf{vec}\left(\mathcal{Z}_i\right)^{\top} \left(\mathbf{K}_3 \otimes \mathbf{K}_2 \otimes \mathbf{K}_1\right) \mathbf{vec}\left(\mathcal{Z}_j\right).$$

The multi-linear kernel defines a similarity metric between pairs of tensor data. We provide an illustration of the multi-linear kernel in Figure 1b. In this model, $\mathbf{K}_1 \in \mathbb{R}^{h \times h}, \mathbf{K}_2 \in \mathbb{R}^{w \times w}, \mathbf{K}_3 \in \mathbb{R}^{C \times C}$ capture the mode-specific covariance structure.

Combining (3), (4) and (5) together, our method essentially specifies the following tensor GP with a new kernel $\mathcal{K}(\cdot, \cdot)$:

$$y_i = f(\mathcal{X}_i) + \epsilon_i, \quad f(\cdot) \sim GP(0, \mathcal{K}(\cdot, \cdot))$$
 (6)

$$\mathcal{K}(\mathcal{X}_{i}, \mathcal{X}_{i}) = \mathbf{vec}\left(\mathcal{X}_{i}\right)^{\top} \left(\mathbf{K}_{3} \otimes \mathbf{K}_{2}^{*} \otimes \mathbf{K}_{1}^{*}\right) \mathbf{vec}\left(\mathcal{X}_{i}\right) \tag{7}$$

$$\mathbf{K}_2^* = \mathbf{B}^{\top} \mathbf{K}_2 \mathbf{B}, \quad \mathbf{K}_1^* = \mathbf{A}^{\top} \mathbf{K}_1 \mathbf{A}, \tag{8}$$

and we call the framework **Tensor G**aussian **Process** with **Spatial Transformation (Tensor-GPST)**, where **A** and **B** transform, in a bi-linear way, the spatial information contained in the imaging data.

Another way of expressing the model is via tensor regression (1) on the original tensor \mathcal{X} . Equivalently, we assume a Gaussian prior over \mathcal{W} :

$$\mathbf{vec}(\mathcal{W}) \sim (\mathbf{I}_C \otimes \mathbf{B} \otimes \mathbf{A})^{\top} \mathbf{vec}(\mathcal{T}), \mathcal{T} \in \mathbb{R}^{h \times w \times C},$$
$$\mathbf{vec}(\mathcal{T}) \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_3 \otimes \mathbf{K}_2 \otimes \mathbf{K}_1), \tag{9}$$

which is similar to a tensor factor model (Chen et al., 2020) coupled with a Gaussian factor with Kronecker-product covariance structure.

2.2. Estimating Algorithm

To estimate the model parameters of Tensor-GPST in (6)-(8), including the tensor contracting factors (\mathbf{A}, \mathbf{B}) , the multi-linear kernel factors $(\mathbf{K}_1, \mathbf{K}_2, \mathbf{K}_3)$, and the idiosyn-

cratic noise variance σ^2 , we minimize the negative marginal Gaussian log-likelihood $\ell(\mathbf{y}|\mathbf{A}, \mathbf{B}, \mathbf{K}_1, \mathbf{K}_2, \mathbf{K}_3, \sigma)$:

$$\ell = \frac{1}{2} \log |\mathbf{K} + \mathbf{D}_{\sigma}| + \frac{1}{2} \mathbf{y}^{\top} (\mathbf{K} + \mathbf{D}_{\sigma})^{-1} \mathbf{y} + \text{const.}, (10)$$

where **K** is an $N \times N$ empirical kernel gram matrix computed using the kernel function (7) for all pairs of tensor data and $\mathbf{D}_{\sigma} = \sigma^2 \mathbf{I}_N$.

To speed up the computation, we approximate each multilinear kernel factor with a factorized form:

$$\mathbf{K}_1 = \mathbf{U}_1^{\mathsf{T}} \mathbf{U}_1, \mathbf{K}_2 = \mathbf{U}_2^{\mathsf{T}} \mathbf{U}_2, \mathbf{K}_3 = \mathbf{U}_3^{\mathsf{T}} \mathbf{U}_3,$$
 (11)

where $\mathbf{U}_1 \in \mathbb{R}^{r_1 \times h}, \mathbf{U}_2 \in \mathbb{R}^{r_2 \times w}, \mathbf{U}_3 \in \mathbb{R}^{r_3 \times C}$. $\mathbf{U}_1, \mathbf{U}_2, \mathbf{U}_3$ are orthogonal matrices with $r_1 \leq h, r_2 \leq w, r_3 \leq C$. The tuning parameter is set as such that $r_1 = h, r_2 = w, r_3 = C$ throughout the paper but can be set to smaller values to enforce a low-rank constraint. With the factorization assumption, one can decompose the gram matrix \mathbf{K} as $\widetilde{\mathbf{U}}\widetilde{\mathbf{U}}^{\top}$, where:

$$\widetilde{\mathbf{U}} = \widetilde{\mathcal{X}}^{\top} \left(\mathbf{I}_C \otimes \mathbf{B} \otimes \mathbf{A} \right)^{\top} \left(\mathbf{U}_3 \otimes \mathbf{U}_2 \otimes \mathbf{U}_1 \right)^{\top},$$

where $\widetilde{\mathcal{X}} = [\mathbf{vec}(\mathcal{X}_1); \mathbf{vec}(\mathcal{X}_2); \dots; \mathbf{vec}(\mathcal{X}_N)]$. The factorized form of \mathbf{K} can simplify the computation of the gradients since one can invert the covariance matrix $(\mathbf{K} + \mathbf{D}_{\sigma})$ with the Woodbury identity, as shown in Appendix B. The computational complexity of the algorithm is thus reduced from the canonical $\mathcal{O}(N^3)$ to $\mathcal{O}(N^2D)$, where D = HWC is the dimension of the data tensor.

Since the tensor contracting factors (\mathbf{A}, \mathbf{B}) are extracting spatial features from each channel of \mathcal{X}_i , we assume that each spatial feature can be constructed from several spatially-contiguous regions for better interpretability. This leads us to the assumption that each feature map $\mathbf{W}_{s,t} = \boldsymbol{\alpha}_s^{\top} \boldsymbol{\beta}_t$ has certain degrees of spatial smoothness. We introduce the spatial smoothness assumption into our model via regularizing its anisotropic total variation norm $\|\mathbf{W}_{s,t}\|_{\text{TV}}$:

$$\|\mathbf{W}_{s,t}\|_{\text{TV}} = \sum_{i=1}^{H-1} \sum_{j=1}^{W} |\mathbf{W}_{s,t}(i+1,j) - \mathbf{W}_{s,t}(i,j)| + \sum_{i=1}^{H} \sum_{j=1}^{W-1} |\mathbf{W}_{s,t}(i,j+1) - \mathbf{W}_{s,t}(i,j)|.$$

A more general class of total variation norm penalty on tensor regression model coefficients can be found in Wang et al. (2017). In Lemma 2.1, we derive a simplified form of $\|\mathbf{W}_{s,t}\|_{\mathrm{TV}}$, making the estimation of \mathbf{A} and \mathbf{B} easier.

Lemma 2.1. The anisotropic total variation (TV) norm on feature map $\{\mathbf{W}_{s,t}\}_{s=1,t=1}^{h,w}$ induces a fused-lasso (Tibshi-

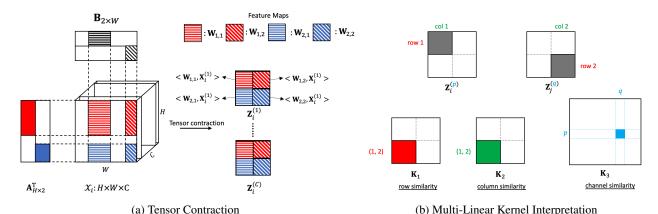


Figure 1: (a) Example of the tensor contraction step for tensor data $\mathcal{X}_i \in \mathbb{R}^{H \times W \times C}$ to its latent tensor $\mathcal{Z}_i \in \mathbb{R}^{2 \times 2 \times C}$. The tensor contracting factors \mathbf{A} , \mathbf{B} are sparse (colored/dashed bands indicate nonzero elements) and they jointly extract features from $\mathbf{X}_i^{(1)}, \dots, \mathbf{X}_i^{(C)}$ with rank-1 feature maps $\{\mathbf{W}_{1,1}, \mathbf{W}_{1,2}, \mathbf{W}_{2,1}, \mathbf{W}_{2,2}\}$. Each channel of $\mathbf{Z}_i^{(c)}$ has 2×2 features, based on the inner product of every feature map with the channel data $\mathbf{X}_i^{(c)}$. (b) Example of the multi-linear kernel with a pair of latent tensor data $(\mathcal{Z}_i, \mathcal{Z}_j)$. Any pair of pixels in \mathcal{Z}_i and \mathcal{Z}_j , e.g. $\mathbf{Z}_i^{(p)}(1,1)$ and $\mathbf{Z}_j^{(q)}(2,2)$ in the plot (colored in gray), are weighted by the product of their row similarity $\mathbf{K}_1(1,1)$ (red), column similarity $\mathbf{K}_2(2,2)$ (green) and channel similarity $\mathbf{K}_3(p,q)$ (blue), in the kernel function (5) for defining the similarity of $\mathcal{Z}_i, \mathcal{Z}_j$. See (44) for a formulaic explanation.

rani et al., 2005) penalization on A (and B), namely:

$$\sum_{s=1}^{h} \sum_{t=1}^{w} \|\mathbf{W}_{s,t}\|_{\text{TV}} = \|\nabla_x \mathbf{B}\|_1 \|\mathbf{A}\|_1 + \|\mathbf{B}\|_1 \|\nabla_x \mathbf{A}\|_1,$$
(12)

where ∇_x computes the horizontal gradient of a matrix, i.e. $\nabla_x \mathbf{A}_{m \times n}(i,j) = \mathbf{I}_{\{j \neq n\}} [\mathbf{A}(i,j+1) - \mathbf{A}(i,j)]$, and $\|\cdot\|_1$ is the elementwise 1-norm of a matrix.

We leave the proof to Appendix A.

The fused-lasso penalty penalizes the sparsity and smoothness of **A**, weighted by the smoothness and sparsity of **B** and vice versa. Jointly, our estimating problem is attempting to minimize the following penalized negative log-likelihood:

$$\begin{split} L(\mathbf{y}|\mathbf{A},\mathbf{B},\mathbf{U}_{1:3},\sigma) &= \ell(\mathbf{y}|\mathbf{A},\mathbf{B},\mathbf{U}_{1:3},\sigma) + \lambda \mathbf{R}(\mathbf{A},\mathbf{B}),\\ \text{where } \mathbf{R}(\mathbf{A},\mathbf{B}) &= \left\|\nabla_x\mathbf{B}\right\|_1 \|\mathbf{A}\|_1 + \|\mathbf{B}\|_1 \left\|\nabla_x\mathbf{A}\right\|_1 \text{ and }\\ \mathbf{U}_{1:3} \text{ is the collection of } \mathbf{U}_1,\mathbf{U}_2,\mathbf{U}_3. \end{split}$$

The total variation penalty can create feature maps with sharp edges and leads to sparsity for interpretation. In the estimating algorithm, we use proximal gradient descent to estimate the tensor contracting factors (\mathbf{A}, \mathbf{B}) and cyclically update the parameters in the order of: $\mathbf{A} \to \mathbf{B} \to (\mathbf{U}_1, \mathbf{U}_2, \mathbf{U}_3) \to \sigma \to \mathbf{A} \to \dots$. The fused-lasso penalty over \mathbf{A} and \mathbf{B} makes the proximal step a well-defined *fused lasso 1-D signal approximation* problem (Friedman et al., 2007). Specifically, at the $(i+1)^{\text{th}}$ iteration, we first propose a gradient descent update for \mathbf{A} , denoted as $\widehat{\mathbf{A}}^{(i+\frac{1}{2})}$,

with stepsize η_i . The final updated value for **A**, i.e. $\widehat{\mathbf{A}}^{(i+1)}$, is the minimizer of the proximal step:

$$\begin{split} \widehat{\mathbf{A}}^{(i+1)} &= \mathbf{prox}_{\mathrm{TV}} \left(\widehat{\mathbf{A}}^{(i+\frac{1}{2})} \right) \\ &= \operatorname*{arg\,min}_{\mathbf{A}} \left\{ \frac{1}{2\eta_i} \left\| \mathbf{A} - \widehat{\mathbf{A}}^{(i+\frac{1}{2})} \right\|_{\mathrm{F}}^2 + \lambda \mathrm{R}(\mathbf{A}, \widehat{\mathbf{B}}^{(i)}) \right\}, \end{split}$$

which can be easily solved by first solving the minimization without the ℓ_1 -penalty on ${\bf A}$ and then apply a soft-thresholding operator to obtain the exact minimizer (see Proposition 1 of Friedman et al. (2007) for the justification). The same procedure applies when one updates ${\bf B}$. We summarize the outline of the estimating algorithm in Algorithm 1 and provide the details of the derivation of gradients and the proximal step in Appendix ${\bf B}$.

Since any pair of (\mathbf{A}, \mathbf{B}) can be re-scaled by a constant c_1 such that: $(\mathbf{B} \otimes \mathbf{A}) = (c_1^{-1}\mathbf{B}) \otimes (c_1\mathbf{A})$, we re-scale the norm of $(\widehat{\mathbf{A}}^{(i)}, \widehat{\mathbf{B}}^{(i)})$ after each iteration to ensure that there is no scaling identifiability issue for the tensor contraction operation.

We do not enforce the orthonormality of $\mathbf{U}_1, \mathbf{U}_2, \mathbf{U}_3$, but a good initialization can still obtain reasonable approximations according to Yu et al. (2018). To give a warm start of the model parameters, one can consider solving a tensor regression problem and a tucker decomposition problem

subsequently, as inspired by (9):

$$\min_{\substack{\mathcal{T} \in \mathbb{R}^{h \times w \times C} \\ \mathbf{A} \in \mathbb{R}^{h \times H} \\ \mathbf{B} \subset \mathbb{R}^{w \times W}}} \sum_{i=1}^{N} \left(y_i - \left\langle \mathcal{X}_i, \mathcal{T} \times_1 \mathbf{A}^\top \times_2 \mathbf{B}^\top \right\rangle \right)^2, \quad (14)$$

$$\min_{\substack{\mathcal{S} \in \mathbb{R}^{r_1 \times r_2 \times r_3} \\ \mathbf{U}_1, \mathbf{U}_2, \mathbf{U}_3}} \left\| \widehat{\mathcal{T}} - \mathcal{S} \times_1 \mathbf{U}_1^\top \times_2 \mathbf{U}_2^\top \times_3 \mathbf{U}_3^\top \right\|^2.$$
(15)

One obtains $\widehat{\mathbf{A}}^{(0)}$, $\widehat{\mathbf{B}}^{(0)}$ from (14) and $\widehat{\mathbf{U}}_{1:3}^{(0)}$ from (15).

Algorithm 1 Alternating Proximal Gradient Descent Algorithm for Tensor-GPST Estimation

Initialize $\widehat{\mathbf{A}}^{(0)}, \widehat{\mathbf{B}}^{(0)}, \widehat{\mathbf{U}}_1^{(0)}, \widehat{\mathbf{U}}_2^{(0)}, \widehat{\mathbf{U}}_3^{(0)}, \widehat{\sigma}^{(0)}$ randomly. Set iteration counter $i \leftarrow 0$. while not converge and $i \leq \max$ -iter do $\widehat{\mathbf{A}}^{(i+\frac{1}{2})} \leftarrow \widehat{\mathbf{A}}^{(i)} - \eta_i \nabla_{\mathbf{A}} \ell(\mathbf{y} | \widehat{\mathbf{A}}^{(i)}, \widehat{\mathbf{B}}^{(i)}, \widehat{\mathbf{U}}_{1:3}^{(i)}, \widehat{\sigma}^{(i)}).$ $\widehat{\mathbf{A}}^{(i+1)} \leftarrow \mathbf{prox}_{\mathsf{TV}}(\widehat{\mathbf{A}}^{(i+\frac{1}{2})}).$ // Fused-Lasso $\widehat{\mathbf{B}}^{(i+\frac{1}{2})} \leftarrow \widehat{\mathbf{B}}^{(i)} - \eta_i \nabla_{\mathbf{B}} \ell(\mathbf{y} | \widehat{\mathbf{A}}^{(i+1)}, \widehat{\mathbf{B}}^{(i)}, \widehat{\mathbf{U}}_{1:3}^{(i)}, \widehat{\sigma}^{(i)}).$ $\widehat{\mathbf{B}}^{(i+1)} \leftarrow \mathbf{prox}_{\mathsf{TV}}(\widehat{\mathbf{B}}^{(i+\frac{1}{2})}).$ // Fused-Lasso Re-scale $\widehat{\mathbf{A}}^{(i+1)}, \widehat{\mathbf{B}}^{(i+1)}$ s.t. $\|\widehat{\mathbf{A}}^{(i+1)}\|_{\mathsf{F}} = 1.$ for $\mathbf{j}=1:3$ do $\mathbf{G}_j \leftarrow \nabla_{\mathbf{U}_j} \ell(\mathbf{y} | \widehat{\mathbf{A}}^{(i+1)}, \widehat{\mathbf{B}}^{(i+1)}, \widehat{\mathbf{U}}_{-j}^{(i)}, \widehat{\sigma}^{(i)}).^1$ $\widehat{\mathbf{U}}_j^{(i+1)} \leftarrow \widehat{\mathbf{U}}_j^{(i)} - \eta_i \mathbf{G}_j.$ end for $t \leftarrow \nabla_{\sigma} \ell(\mathbf{y} | \widehat{\mathbf{A}}^{(i+1)}, \widehat{\mathbf{B}}^{(i+1)}, \widehat{\mathbf{U}}_{1:3}^{(i+1)}, \widehat{\sigma}^{(i)}).$ $\widehat{\sigma}^{(i+1)} \leftarrow \widehat{\sigma}^{(i)} - \eta_i t.$ $i \leftarrow i+1$ end while Output: $\widehat{\mathbf{A}}^{(i)}, \widehat{\mathbf{B}}^{(i)}, \widehat{\mathbf{U}}_{1:3}^{(i)}, \widehat{\sigma}^{(i)}$

2.3. Convergence Analysis

In this subsection, we provide the convergence analysis of Algorithm 1. Theorem 2.2 provides the upper bound of the loss function (13), evaluated at the estimators output by the algorithm, with respect to its global minimum. We show that the total variation penalty and the alternating proximal gradient descent introduce extra gaps between the achieved loss and its global minimum.

Theorem 2.2. Given the loss function $L(\cdot)$ in (13), assume that the negative log-likelihood $\ell(\cdot)$ is convex for any of the four parameter blocks: $\{A\}, \{B\}, \{U_{1:3}\}, \{\sigma\},$ with the other three blocks being fixed, and the gradients of $\ell(\cdot)$ are Lipschitz continuous with Lipschitz constant: M_A, M_B, M_U, M_σ , respectively. Then with a constant learning rate $\alpha \leq 1/\max\{M_A, M_B, M_U, M_\sigma\}$, the alternating proximal gradient descent algorithm in Algorithm

1 leads to the following upper bound on the loss function $L(\cdot)$:

$$4(K+1)\left[L(\widehat{\boldsymbol{\theta}}^{(K+1)}) - L(\boldsymbol{\theta}^*)\right] \leq \frac{\delta^{(0)}}{2\alpha}$$

$$+ \sum_{k=0}^{K} h_{\lambda}(\widehat{\mathbf{A}}^{(k+1)} - \mathbf{A}^*, \widehat{\mathbf{B}}^{(k+1)} - \mathbf{B}^*, \widehat{\mathbf{B}}^{(k)} - \mathbf{B}^*)$$

$$+ \frac{1}{2\alpha} \sum_{k=0}^{K} \tau\left(\widehat{\boldsymbol{\theta}}^{(k+1)}, \widehat{\boldsymbol{\theta}}^{(k)}, \boldsymbol{\theta}^*\right), \tag{16}$$

where $\widehat{\boldsymbol{\theta}}^{(k)} = \left\{ \widehat{\mathbf{A}}^{(k)}, \widehat{\mathbf{B}}^{(k)}, \widehat{\mathbf{U}}^{(k)}_{1:3}, \widehat{\sigma}^{(k)} \right\}$ and $\boldsymbol{\theta}^*$ is the global minimizer of $L(\cdot)$. $\delta^{(0)}$ is the squared ℓ_2 initialization error, $h_{\lambda}(\cdot) \geq 0$ is the total-variation gap (TV-gap) and $\tau(\cdot) \geq 0$ is the alternating descent gap (ALT-gap). K is the total number of iterations. As a result, if one has any three blocks of parameters fixed at their global minima, the remaining block will converge to its global minima at the rate of $\mathcal{O}(1/K)$, which echoes the convergence rate of (proximal) gradient descent.

We leave the proof to Appendix C and make a few remarks. Remark 2.3. As $\widehat{\mathbf{A}}^{(k)} \to \mathbf{A}^*$, $\widehat{\mathbf{B}}^{(k)} \to \mathbf{B}^*$, one has $h_{\lambda}(\cdot) \to 0$. The TV-gap is incurred because we alternatively update \mathbf{A} and \mathbf{B} , and using the current iteration's estimate of \mathbf{A} (or \mathbf{B}) for updating \mathbf{B} (or \mathbf{A}) with the total variation penalty leads to extra errors. See the definition of $h_{\lambda}(\cdot)$ in (41).

Remark 2.4. As $\widehat{\boldsymbol{\theta}}^{(k)} \to \boldsymbol{\theta}^*$, $\tau(\cdot) \to 0$. The ALT-gap $\tau(\cdot)$ arises because we use the current iteration's estimate for all but one block of parameters to estimate the gradient of the block of interest. If the algorithm terminates at a local minima, the non-vanishing TV-gap and ALT-gap leads to a non-zero gap for the achieved loss from the global minimum. See the definition of $\tau(\cdot)$ in (42).

Remark 2.5. Tensor regression models with Tucker-type low-rankness have non-convex negative-likelihood function $\ell(\cdot)$ (Li et al., 2018). But conditioning on all but one block of parameter, $\ell(\cdot)$ is convex for each individual block. We do not verify the convexity of $\ell(\cdot)$ in our particular model due to the complexity of the kernel function. Empirically, as we show in Figure 5 in Appendix C and also in Yu et al. (2018), such alternating gradient descent algorithm works well with the optimization problem and the loss function decays at the rate of $\mathcal{O}(1/K)$.

3. Experiments

In this section, we validate our method via both simulation studies and an application to an astrophysics dataset for solar flare forecasting. We also compare our method against other benchmark tensor regression models. In particular, we are interested in applications to imaging data where the

 $^{^{1}}$ We use \mathbf{U}_{-j} to denote the collection of $\mathbf{U}_{1},\mathbf{U}_{2},\mathbf{U}_{3}$ but exclude $\mathbf{U}_{j}.$

predictive signals appear in different channels and various locations within a channel. Such patterns are common in astrophysical imaging data where the solar flare precursors could appear in the images collected by the astrophysical telescopes at various frequencies and the locations of the precursors could be random within an image channel.

3.1. Simulation Study

We simulate a tensor dataset $\{\mathcal{X}_i\}_{i=1}^N$ with each \mathcal{X}_i having 3 channels of size 25×25 . For each $25 \times 25 \times 3$ tensor data \mathcal{X}_i , we randomly pick one of the three channels as the *signal* channel, with equal probability, and the remaining two channels as the *noise* channels. The noise channel contains i.i.d. pixels from $\mathcal{N}(0,0.3)$, and the signal channel uses the same background noise distribution except having a 5×5 signal block that contains i.i.d. pixels from $\mathcal{N}(4,0.3)$. The location of the 5×5 block is fixed at the center of the 25×25 image if channel 2 is the signal channel (see Type 2 in Figure 2), and is randomly picked at one of the four corners (top-left, top-right, bottom-left, bottom-right) if channel 1 or 3 is the signal channel (see Type 1 and 3 in Figure 2).

We simulate the tensor contracting factors $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{3 \times 25}$ with a banded structure, leading to a $3 \times 3 \times 3$ contracted tensor \mathcal{Z} , such that \mathbf{A} and \mathbf{B} are extracting features from the 5×5 blocks with signal, see the bottom of Figure 2 for the example of the contracted tensor \mathcal{Z} . The multi-linear kernel setup and the generating process of the regression labels $\{y_i\}_{i=1}^N$ are detailed in Appendix D. Generally, channel 2 is simulated such that it is negatively correlated with channels 1 & 3, and channels 1 & 3 are nearly perfectly correlated. As a result, Type 1 & 3 tensors have similar regression labels and differ from those of Type 2.

With the simulation setups, we compare our model against these baseline tensor regression models: Tensor-GP (GP) (Yu et al., 2018), ℓ_2 -regularized tensor regression with CAN-DECOMP/PARAFAC (CP) tensor rank constraints (Guo et al., 2011) (CP), ℓ_2 -regularized tensor regression with coefficient tensor following a Tucker decomposition (Li et al., 2018) (Tucker). In order to check the sensitivity of the choice of the kernel, we also fit the GP model to the vectorized tensor data with a squared-exponential kernel (SE). To showcase the effectiveness of tensor contraction, we also consider fitting a model with a tensor contraction step followed by a GP with squared-exponential kernel for the vectorized, reduced-sized tensor (SE+TC). For simplicity, we implement **SE+TC** by training the model, without the total variation penalty, in an end-to-end fashion with the GPvTorch and Tensorly-Torch packages in Python. Both models involving the SE kernel have automatic relevance determination (ARD) length scales (Bishop & Nasrabadi, 2006).

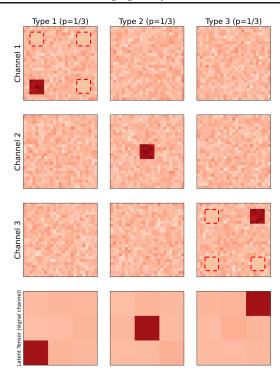


Figure 2: Three types of the simulated tensor data $(\mathcal{X}_i \in \mathbb{R}^{25 \times 25 \times 3})$. Each column is a type (Type 1,2,3) and every sample has equal probability of being one of the three types. Each row (row 1-3) is a data channel (channel 1,2,3). Type 1, 2 and 3 have their *signal* channel in channel 1, 2 and 3, respectively. But the location of the 5×5 signal block is positioned differently. Type 2 has the signal fixed at the center, while type 1 and 3 has the signal placed, with equal probability, in one of the four corners (dashed block shows the other three possible locations). Samples shown are one realization of the simulation. The latent tensor \mathcal{Z} 's signal channel is shown at the bottom. See details in Appendix D.

We simulate the data with size $N \in \{200, 500\}$ and use 75% for training and 25% for testing and compare the rooted-mean-squared-error (RMSE) on both training and testing across all models above as well as our own Tensor-GPST model (GPST). We set the latent tensor dimension as $3\times 3\times 3$ for GPST and SE+TC and the rank for $\mathbf{K}_1,\mathbf{K}_2,\mathbf{K}_3$ of GP as 3 and the CP rank as 9 for CP and the multi-linear rank as $3\times 3\times 3$ for Tucker such that the low-rankness is comparable across all methods. We select the regularization tuning parameter for all models with hyperparameters by 5-fold cross validation. The simulation experiment is iterated 10 times and the testing RMSE is shown in Table 1.

The Tensor-GP (**GP**) method has relatively worse performance on the testing set compared to other low-rank tensor regression methods such as **CP** and **Tucker**. Our method, namely **GPST**, achieves similar performance to the low-

Table 1: Test prediction RMSE for simulated Data for various tensor regression models. 95% confidence interval after \pm . Results are based on 10 repeated runs.

Model	N = 200	N = 500		
GP	0.728 ± 0.125	0.664 ± 0.131		
CP	0.550 ± 0.100	0.548 ± 0.054		
Tucker	0.589 ± 0.206	0.568 ± 0.107		
SE	2.504 ± 2.672	3.275±4.204		
SE+TC	0.627 ± 0.169	0.587 ± 0.098		
GPST (Our Method)	0.578 ± 0.107	0.552 ± 0.076		

rank tensor regression methods (not statistically significantly worse). The GP with vectorized tensor data and squaredexponential kernel, namely SE, performs extremely poorly, which reveals the fact that by vectorizing tensor data one loses the essential structural information of the data. This result necessitates the choice of kernel that is suitable for tensor data, such as the tensor GP. After adding an extra tensor contraction step, the GP with squared-exponential kernel (i.e. SE+TC) performs relatively close to the lowrank tensor regression methods as well as our GPST and is better than the tensor GP. This further suggests that regardless of the kernel choice, the tensor contraction step can boost the performance of GP regression models with tensor covariates. Effectively, the tensor contraction step extracts useful features from the original tensor data for regression, so even if one vectorizes the reduced-sized tensor, one does not lose as much information as the case where tensor contraction is not being used. Finally, we note that with a large sample size (N = 500), the prediction RMSE of the test set gets smaller for all methods but SE.

To make further comparisons of the variants of different Gaussian Process models listed in Table 1 on their ability to quantify the uncertainties of the predictions made, we compare these GP models' mean standardized log loss (MSLL) (Williams & Rasmussen, 2006), as defined below:

$$MSLL = \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} \left\{ \frac{1}{2} \log \left(2\pi \hat{\sigma}^2 \right) + \frac{(y_i - \hat{y}_i)^2}{2\hat{\sigma}^2} \right\},$$

where \hat{y}_i is the predicted label for the i^{th} testing sample and $\hat{\sigma}$ is the estimated standard deviation of the noise term. Generally speaking, a smaller MSLL indicates a better testing prediction. We list the testing set MSLL for **GP**, **SE**, **SE+TC** and **GPST** in Table 2.

The result reveals that our method has statistically significantly smaller MSLL, as indicated by a one-sided paired t-test, compared to the other methods under both sample sizes. Also, the models with tensor contraction, including **SE+TC** and **GPST**, have smaller MSLL compared to their counterparts without tensor contraction, which further sug-

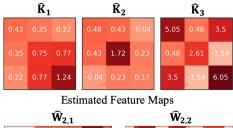
Table 2: Test Mean Standardized Log Loss (MSLL) for the 4 variants of GP models. 95% confidence interval after \pm . Results are based on 10 repeated runs.

Model	N = 200	N = 500
GP	1.162 ± 0.439	1.092 ± 0.421
SE	$2.457{\pm}1.898$	2.999 ± 3.167
SE+TC	0.972 ± 0.214	0.919 ± 0.123
GPST	$\boldsymbol{0.882 \pm 0.201}$	$\boldsymbol{0.835 \pm 0.156}$

gests that tensor contraction can be helpful for reducing the errors made by GP models with tensor data.

The estimators of the multi-linear kernel factors $\mathbf{K}_1, \mathbf{K}_2, \mathbf{K}_3$ and the feature maps of the Tensor-GPST model with $\lambda=1.0$ for one random simulation dataset are visualized in Figure 3. One can see that the feature map $\widehat{\mathbf{W}}_{2,2}$ and $\widehat{\mathbf{W}}_{3,2}$ capture the corner and center blocks, and the covariances between the two feature maps are also high, as suggested by $\widehat{\mathbf{K}}_1(2,3)=0.77$ and $\widehat{\mathbf{K}}_2(2,2)=1.72$. Channels 1 & 3 have high covariances $(\widehat{\mathbf{K}}_3(1,3)=3.5)$, indicating that they share similar "corner signal" patterns and coincides with our ground truth setup (see Figure 6a for the ground truth of \mathbf{K}_3).

Estimated Multi-Linear Kernel



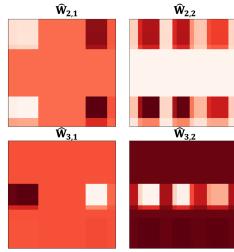


Figure 3: Estimated kernels (top) and non-zero feature maps (bottom) by **GPST** with $\lambda=1.0$ for one random simulation dataset.

Overall, the simulation experiments convey two messages:

- Adding the tensor contraction step leads to better regression performances when the signals have low-rank structures, robust to the choice of kernel, and the performance is similar to other low-rank tensor regression methods such as CP and Tucker.
- The anisotropic total variation penalty, though may not fully recover the underlying sparsity of the tensor contracting factors, can improve the regression performance of Tensor-GPST and also provides more direct interpretations.

The inferior performance of Tensor-GP (**GP**), however, is not suggesting that it is an inferior version of GP when dealing with tensor data, as we have demonstrated by comparing it against the GP with squared-exponential kernel (**SE**). The simulation pattern in Figure 2 contains randomness of the signal, making it more beneficial to extract features first using feature maps that cover multiple areas. Directly modeling the covariance structures among all pixels can be difficult in such scenarios.

3.2. Application to Solar Flare Forecasting

A solar flare is an intense localized eruption of electromagnetic radiation in the Sun's atmosphere. Solar flares with high-energy radiation emission can strongly impact the Earth's space weather and potentially interfere the radio communication of the Earth. Recent works on solar flare forecasting (Bobra & Couvidat, 2015; Chen et al., 2019; Wang et al., 2020; Jiao et al., 2020; Sun et al., 2022) have demonstrated the effectiveness of using machine learning algorithms for forecasting flares, using either multivariate time-series data in the form of physical parameters or imaging data provided by the Solar Dynamics Observatory (SDO)/Helioseismic and Magnetic Imager (HMI) (Scherrer et al., 2012) and SDO/Atmospheric Imaging Assembly (AIA) (Lemen et al., 2012). It has been shown that these imaging data have low-dimensional representations that contain flare discriminating signals (Sun et al., 2021). Our methodology makes the astrophysical interpretation more accessible as compared to the previous deep learning approaches in that our model has a much shallower structure with only a feature extraction layer and a regression layer.

Here, we consider the specific problem of forecasting the intensity of a solar flare. In our dataset, we have 1, 329 flare samples from year 2010 to 2018, consisting of a total of 479 M-class and X-class flares and 850 B-class flares. The class of a flare is determined by the X-ray peak brightness in the range of 1-8Å. The B-class flare has its brightness within $10^{-7} \sim 10^{-6} \text{W/m}^2$, which is considered weak and barely harmful, while the minimum M-class and X-class

flares have brightness above 10^{-5} and 10^{-4}W/m^2 , respectively. These more energetic flares are capable of heating and ionizing the the upper atmosphere, resulting in brief radio blackouts and increased satellite drag. We collect the AIA-HMI imaging data for each flare, 1 hour prior to its peak, and each flare data is a 10-channel tensor data of size $50 \times 50 \times 10$, where spatial dimensions are binned down by roughly a factor of 10. We leave the data preprocessing steps and the astrophysical background to Appendix E.

Our goal here is to utilize the 10-channel tensor data \mathcal{X}_i to predict the flare intensity y_i and find the discriminating factors for M/X-class and B-class flares. We randomly split our dataset into a 75% training set (359 M/X/637 B) and a 25% testing set (120 M/X/213 B), and centering after log-transforming the flare intensity such that the B-class flare has $y_i \leq -0.5$ and M/X-class flare has $y_i \geq 0.5$.

We report the solar flare intensity prediction result across four different models: Tensor-GP (GP), Tensor-GPST (**GPST**), tensor regression with CP rank constraints (**CP**) and tensor regression with Tucker decomposition form (Tucker). The hyperparameters are set such that the models have the same latent dimensionality $(3 \times 3 \times 3)$ or the rank (9 for **CP** and $3 \times 3 \times 3$ for **Tucker**) of the regression coefficients. The metrics used are rooted mean-squared error (RMSE), R-squared and MSLL. Additionally, we consider transforming the regression model to a binary classification model by thresholding the prediction at 0.0 such that any $\hat{y}_i \geq 0$ indicates an M/X-class flare and any $\hat{y}_i < 0$ indicates a B-class flare. Then we evaluate the resulting binary classification model with the True Skill Statistics (TSS)². A skillful binary classifier for weak vs. strong solar flare is desirable for operational use. Results on the training and testing set are summarized in Table 3, with 10 random train/test splits.

Tensor-GP (**GP**) shows worse generalizability on the testing data as compared to the other three methods. **GPST** has slightly better testing set performance compared to **CP** and **Tucker**, but is not statistically significantly better than **Tucker**. Similar to the simulation data, the flare data exhibits randomness of the location of flare predictive signals, making the tensor contraction a critical step for improving the Tensor-GP method.

In Figure 4, we visualize the class-average AIA-131Å in the left column. There is a stark contrast between the two flare classes for this channel and many other channels as we show in Appendix F. A convenient output of our Tensor-GPST model is the direct estimation of channel covariances in the multi-linear kernel, and we visualize the estimated $\widehat{\mathbf{K}}_3$ in the Figure as well. The estimated $\widehat{\mathbf{K}}_3$ reveals the important

²True Skill Statistics is defined as: TSS = TP/(TP+FN) – FP/(FP+TN), where TP, TN, FP, FN represents true positive, true negative, false positive and false negative in the confusion matrix.

Table 3: Solar flare intensity regression performance on the training and testing sets for four tensor regression models. Results based on 10 random splits and 95% confidence intervals are provided after \pm .

Model	Training (75% of the samples)			Testing (25% of the samples)				
	RMSE	\mathbb{R}^2	MSLL	TSS	RMSE	\mathbb{R}^2	MSLL	TSS
GP	0.646±0.019	0.336 ± 0.044	1.028 ± 0.134	0.466 ± 0.039	0.772±0.239	0.182 ± 0.114	1.138 ± 0.085	0.362±0.159
CP	0.564 ± 0.035	0.501 ± 0.077	_	$\boldsymbol{0.625 \pm 0.069}$	0.706 ± 0.051	0.230 ± 0.078	_	0.398 ± 0.092
Tucker	0.679 ± 0.014	0.269 ± 0.028	_	$0.426{\pm}0.052$	0.683 ± 0.040	0.259 ± 0.079	_	0.414 ± 0.134
GPST	0.661 ± 0.014	0.305 ± 0.023	1.005 ± 0.021	0.449 ± 0.040	0.681±0.043	$0.265{\pm}0.087$	$\boldsymbol{1.035 \pm 0.061}$	0.412 ± 0.112

channel pairs when defining the similarity of pairs of tensor data, and we formalize this channel pair importance notion in (45) of Appendix F. To the best of our knowledge, our model is the first to consider the channel interactions for solar flare forecasting.

In the lower right panel of Figure 4, we visualize the pixels that have at least one feature map with weight $> 5 \times 10^{-3}$. These pixels contribute significantly to building the latent tensor, and are thus being considered as the most relevant pixels for solar flare prediction. As one can see, the selected pixels are concentrated around the two brightest spots of the AIA-131Å for the M/X-class and also around the boundary. These pixels contain two most significant flare discriminating factors: 1) the brightest spots of the AIA images; 2) the span of the bright regions (as M/X flares still have large AIA image intensities near the boundary but not B flares).

4. Conclusion

In this paper, we propose a new methodology called **Tensor-GPST** for fitting Gaussian Process Regression (GPR) model on labelled multi-channel imaging data. We propose a tensor contraction operation to reduce the dimensionality of the tensor and also introduce anisotropic total variation penalty to the tensor contraction parameters to allow for interpretable feature extraction. We see improvements on the regression performances over the original Tensor-GP (Yu et al., 2018) in both simulation and the solar flare forecasting task. The capability of the model in generating an interpretable low-dimensional tensor representation makes it ideal for many other scientific applications, such as predicting ADHD with Brain-Image (Li et al., 2018) and studying the association of brain connectivity with human traits (Papadogeorgou et al., 2021).

The current model has several limitations that can potentially lead to future research directions. First, we do not impose explicit identifiability constraints for the tensor contraction parameters and the multi-linear kernel parameters. This makes the optimization problem unconstrained thus enabling a simple gradient-based algorithm, but makes the parameters not fully identifiable. Second, the model has higher computational complexity as compared to the non-GP tensor regression models due to its GP formulation. A

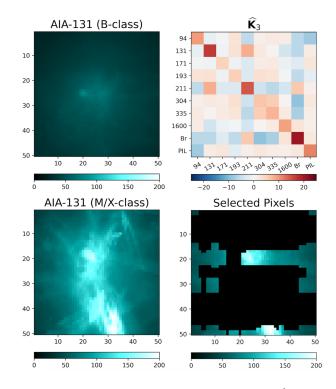


Figure 4: (Left column) The average AIA-131Å for all B-class flares and all M/X-class flares. (Right column) Estimated $\hat{\mathbf{K}}_3$ in the multi-linear kernel that captures the channel-channel covariances (top). Pixels with at least one feature map with weight $> 5 \times 10^{-3}$ (bottom). We visualize the selected pixels with M/X class average AIA-131Å as the background. See full results in Appendix F.

more efficient computational algorithm is needed to handle larger datasets.

Our code is available on GitHub at https://github.com/husun0822/TensorGPST.

Acknowledgement

This material is based upon work supported by the National Science Foundation under Grant NSF-PHY 2027555 and NSF-DMS 113397. WBM is supported by NASA grant 80NSSC18K1208.

References

- Barbero, A. and Sra, S. Modular Proximal Optimization for Multidimensional Total-Variation Regularization. *Journal of Machine Learning Research*, 19(1):2232–2313, 2018.
- Bishop, C. M. and Nasrabadi, N. M. *Pattern recognition* and machine learning, volume 4. Springer, 2006.
- Bobra, M. G. and Couvidat, S. Solar Flare Prediction using SDO/HMI Vector Magnetic Field Data with a Machine-Learning Algorithm. *The Astrophysical Journal*, 798(2): 135, 2015.
- Chen, E. Y., Xia, D., Cai, C., and Fan, J. Semiparametric Tensor Factor Analysis by Iteratively Projected SVD. *arXiv preprint arXiv:2007.02404*, 2020.
- Chen, Y., Manchester, W. B., Hero, A. O., Toth, G., Du-Fumier, B., Zhou, T., Wang, X., Zhu, H., Sun, Z., and Gombosi, T. I. Identifying Solar Flare Precursors using Time Series of SDO/HMI Images and SHARP Parameters. *Space Weather*, 17(10):1404–1426, 2019.
- Friedman, J., Hastie, T., Höfling, H., and Tibshirani, R. Pathwise Coordinate Optimization. *The Annals of Applied Statistics*, 1(2):302–332, 2007.
- Guo, W., Kotsia, I., and Patras, I. Tensor Learning for Regression. *IEEE Transactions on Image Processing*, 21 (2):816–827, 2011.
- Hao, B., Wang, B., Wang, P., Zhang, J., Yang, J., and Sun, W. W. Sparse Tensor Additive Regression. *Journal of machine learning research*, 22, 2021.
- Hung, H. and Wang, C.-C. Matrix Variate Logistic Regression Model with Application to EEG Data. *Biostatistics*, 14(1):189–202, 2013.
- Jiao, Z., Sun, H., Wang, X., Manchester, W., Gombosi, T., Hero, A., and Chen, Y. Solar Flare Intensity Prediction with Machine Learning Models. *Space Weather*, 18(7): e2020SW002440, 2020.
- Jiménez, Á. B. and Sra, S. Fast Newton-Type Methods for Total Variation Regularization. In *International Confer*ence on Machine Learning, 2011.
- Kang, J., Reich, B. J., and Staicu, A.-M. Scalar-on-Image Regression via the Soft-thresholded Gaussian Process. *Biometrika*, 105(1):165–184, 2018.
- Kolda, T. G. and Bader, B. W. Tensor Decompositions and Applications. *SIAM review*, 51(3):455–500, 2009.

- Kossaifi, J., Lipton, Z. C., Kolbeinsson, A., Khanna, A., Furlanello, T., and Anandkumar, A. Tensor Regression Networks. *Journal of Machine Learning Research*, 21(1): 4862–4882, 2020.
- Lemen, J. R., Title, A. M., Akin, D. J., Boerner, P. F., Chou, C., Drake, J. F., Duncan, D. W., Edwards, C. G., Friedlaender, F. M., Heyman, G. F., Hurlburt, N. E., Katz, N. L., Kushner, G. D., Levay, M., Lindgren, R. W., Mathur, D. P., McFeaters, E. L., Mitchell, S., Rehse, R. A., Schrijver, C. J., Springer, L. A., Stern, R. A., Tarbell, T. D., Wuelser, J.-P., Wolfson, C. J., Yanari, C., Bookbinder, J. A., Cheimets, P. N., Caldwell, D., Deluca, E. E., Gates, R., Golub, L., Park, S., Podgorski, W. A., Bush, R. I., Scherrer, P. H., Gummin, M. A., Smith, P., Auker, G., Jerram, P., Pool, P., Soufli, R., Windt, D. L., Beardsley, S., Clapp, M., Lang, J., and Waltham, N. The Atmospheric Imaging Assembly (AIA) on the Solar Dynamics Observatory (SDO). *Solar Physics*, 275(1-2):17–40, January 2012. doi: 10.1007/s11207-011-9776-8.
- Li, X., Xu, D., Zhou, H., and Li, L. Tucker Tensor Regression and Neuroimaging Analysis. *Statistics in Biosciences*, 10(3):520–545, 2018.
- Lock, E. F., Hoadley, K. A., Marron, J. S., and Nobel, A. B. Joint and Individual Variation Explained (jive) for Integrated Analysis of Multiple Data Types. *The Annals of Applied Statistics*, 7(1):523, 2013.
- Papadogeorgou, G., Zhang, Z., and Dunson, D. B. Soft Tensor Regression. *Journal of Machine Learning Research*, 22:219–1, 2021.
- Parikh, N., Boyd, S., et al. Proximal algorithms. *Foundations and trends*® *in Optimization*, 1(3):127–239, 2014.
- Rudin, L. I., Osher, S., and Fatemi, E. Nonlinear Total Variation based Noise Removal Algorithms. *Physica D: nonlinear phenomena*, 60(1-4):259–268, 1992.
- Scherrer, P. H., Schou, J., Bush, R. I., Kosovichev, A. G., Bogart, R. S., Hoeksema, J. T., Liu, Y., Duvall, T. L., Zhao, J., Title, A. M., Schrijver, C. J., Tarbell, T. D., and Tomczyk, S. The Helioseismic and Magnetic Imager (HMI) Investigation for the Solar Dynamics Observatory (SDO). *Solar Physics*, 275(1-2):207–227, January 2012. doi: 10.1007/s11207-011-9834-2.
- Schrijver, C. J. A Characteristic Magnetic Field Pattern Associated with All Major Solar Flares and its Use in Flare Forecasting. *The Astrophysical Journal*, 655(2): L117, 2007.
- Sun, H., Manchester IV, W., and Chen, Y. Improved and Interpretable Solar Flare Predictions with Spatial and Topological Features of the Polarity Inversion Line Masked

- Magnetograms. *Space Weather*, 19(12):e2021SW002837, 2021.
- Sun, Z., Bobra, M. G., Wang, X., Wang, Y., Sun, H., Gombosi, T., Chen, Y., and Hero, A. Predicting Solar Flares using CNN and LSTM on Two Solar Cycles of Active Region Data. *The Astrophysical Journal*, 931(2):163, 2022.
- Tibshirani, R., Saunders, M., Rosset, S., Zhu, J., and Knight, K. Sparsity and Smoothness via the Fused Lasso. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(1):91–108, 2005.
- Tibshirani, R. J. Adaptive Piecewise Polynomial Estimation via Trend Filtering. *The Annals of Statistics*, 42(1):285–323, 2014.
- Wang, X., Zhu, H., and Initiative, A. D. N. Generalized Scalar-on-Image Regression Models via Total Variation. *Journal of the American Statistical Association*, 112(519): 1156–1168, 2017.
- Wang, X., Chen, Y., Toth, G., Manchester, W. B., Gombosi,
 T. I., Hero, A. O., Jiao, Z., Sun, H., Jin, M., and Liu,
 Y. Predicting Solar Flares with Machine Learning: Investigating Solar Cycle Dependence. *The Astrophysical Journal*, 895(1):3, 2020.
- Williams, C. K. and Rasmussen, C. E. Gaussian Processes for Machine Learning, volume 2. MIT press Cambridge, MA, 2006.
- Yu, R., Li, G., and Liu, Y. Tensor Regression meets Gaussian Processes. In *International Conference on Artificial Intelligence and Statistics*, pp. 482–490. PMLR, 2018.
- Zhou, H. and Li, L. Regularized Matrix Regression. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(2):463–483, 2014.
- Zhou, H., Li, L., and Zhu, H. Tensor Regression with Applications in Neuroimaging Data Analysis. *Journal of the American Statistical Association*, 108(502):540–552, 2013.

A. Proof of Lemma 2.1

The anisotropic total variation penalty can be simplified as follows, thanks to the rank-1 assumption on the feature map $W_{s,t}$:

$$\|\mathbf{W}_{s,t}\|_{\text{TV}} = \sum_{i=1}^{H-1} \sum_{j=1}^{W} |\mathbf{W}_{s,t}(i+1,j) - \mathbf{W}_{s,t}(i,j)| + \sum_{i=1}^{H} \sum_{j=1}^{W-1} |\mathbf{W}_{s,t}(i,j+1) - \mathbf{W}_{s,t}(i,j)|$$
(17)

$$= \sum_{i=1}^{H-1} \sum_{j=1}^{W} |\mathbf{A}(s,i+1) - \mathbf{A}(s,i)| \cdot |\mathbf{B}(t,j)| + \sum_{i=1}^{H} \sum_{j=1}^{W-1} |\mathbf{B}(t,j+1) - \mathbf{B}(t,j)| \cdot |\mathbf{A}(s,i)|.$$
(18)

As a result, the total variation penalty has an elegant multiplicative formulation:

$$\lambda \sum_{s=1}^{h} \sum_{t=1}^{w} \|\mathbf{W}_{s,t}\|_{\text{TV}} = \left(\lambda \sum_{t=1}^{w} \sum_{j=1}^{W} |\mathbf{B}(t,j)|\right) \cdot \left(\sum_{s=1}^{h} \sum_{i=1}^{H-1} |\mathbf{A}(s,i+1) - \mathbf{A}(s,i)|\right)$$

$$+ \left(\lambda \sum_{t=1}^{w} \sum_{j=1}^{W-1} |\mathbf{B}(t,j+1) - \mathbf{B}(t,j)|\right) \cdot \left(\sum_{s=1}^{h} \sum_{i=1}^{H} |\mathbf{A}(s,i)|\right)$$

$$= \lambda \cdot \|\mathbf{B}\|_{1} \cdot \|\nabla_{x}\mathbf{A}\|_{1} + \lambda \cdot \|\nabla_{x}\mathbf{B}\|_{1} \cdot \|\mathbf{A}\|_{1}, \tag{19}$$

where ∇_x is the horizontal (i.e. row) first-order derivative operator and $\|\cdot\|_1$ is the matrix ℓ_1 norm. (19) turns out to be a fused-lasso type penalty (Tibshirani et al., 2005) with both a penalty on the sparsity of \mathbf{A} and a penalty on the smoothness of each row of \mathbf{A} . This leads to a rank-1 feature map with row smoothness. Different from the fused-lasso penalty, we only introduce one tuning parameter λ instead of λ_1, λ_2 for the sparsity of smoothness of \mathbf{A} separately. Instead, the tuning parameter for the sparsity of smoothness of \mathbf{A} is re-weighted by the smoothness and sparsity of \mathbf{B} and vice versa, according to (19).

B. Proximal Gradient Descent Algorithm for Tensor-GPST

Given the factorization assumption for the tensor multi-linear kernel factors K_1, K_2, K_3 :

$$\mathbf{K}_1 = \mathbf{U}_1^{\mathsf{T}} \mathbf{U}_1, \mathbf{K}_2 = \mathbf{U}_2^{\mathsf{T}} \mathbf{U}_2, \mathbf{K}_3 = \mathbf{U}_3^{\mathsf{T}} \mathbf{U}_3,$$
 (20)

where $\mathbf{U}_1 \in \mathbb{R}^{r_1 \times h}$, $\mathbf{U}_2 \in \mathbb{R}^{r_2 \times w}$, $\mathbf{U}_3 \in \mathbb{R}^{r_3 \times C}$. One can rewrite the penalized negative log-likelihood loss function in (13) as:

$$\min_{\mathbf{A}, \mathbf{B}, \mathbf{U}_{1:3}, \sigma} L = \frac{1}{2} \log \left| \widetilde{\mathbf{U}} \widetilde{\mathbf{U}}^{\top} + \mathbf{D}_{\sigma} \right| + \frac{1}{2} \mathbf{y}^{\top} \left(\widetilde{\mathbf{U}} \widetilde{\mathbf{U}}^{\top} + \mathbf{D}_{\sigma} \right)^{-1} \mathbf{y}$$

$$+ \lambda \sum_{s=1}^{h} \sum_{t=1}^{w} \left\{ \sum_{i=1}^{H-1} \sum_{j=1}^{W} \left| \mathbf{W}_{s,t}(i+1,j) - \mathbf{W}_{s,t}(i,j) \right| + \sum_{i=1}^{H} \sum_{j=1}^{W-1} \left| \mathbf{W}_{s,t}(i,j+1) - \mathbf{W}_{s,t}(i,j) \right| \right\}$$

$$= \ell(\mathbf{A}, \mathbf{B}, \mathbf{U}_{1:3}, \sigma) + \lambda \sum_{s=1}^{h} \sum_{t=1}^{w} \left\| \mathbf{W}_{s,t} \right\|_{TV},$$
(21)

recall that:

$$\mathbf{y} = [y_1, y_2, \dots, y_N]^\top, \quad \widetilde{\mathcal{X}} = [\mathbf{vec}\left(\mathcal{X}_1\right) : \mathbf{vec}\left(\mathcal{X}_2\right) : \dots : \mathbf{vec}\left(\mathcal{X}_N\right)], \quad \widetilde{\mathbf{U}} = \widetilde{\mathcal{X}}^\top \left(\mathbf{I}_C \otimes \mathbf{B} \otimes \mathbf{A}\right)^\top \left(\mathbf{U}_3 \otimes \mathbf{U}_2 \otimes \mathbf{U}_1\right)^\top.$$

We update the model parameters via a block coordinate descent scheme following the order of: $\mathbf{A} \to \mathbf{B} \to (\mathbf{U}_1, \mathbf{U}_2, \mathbf{U}_3) \to \sigma \to \mathbf{A} \to \mathbf{B} \to \dots$

The derivation of the gradients of $\ell(\cdot)$ w.r.t. $A, B, U_{1:3}, \sigma$ have been made trivial thanks to the factorization assumption

(20). For the $(i, j)^{th}$ element of **A**, for instance, we have its partial derivative as:

$$\frac{\partial \ell}{\partial \mathbf{A}}(i,j) = \mathbf{tr} \left[\left(\frac{\partial \ell}{\partial \widetilde{\mathbf{U}}} \right)^{\top} \left(\frac{\partial \widetilde{\mathbf{U}}}{\partial \mathbf{A}(i,j)} \right) \right], \quad \frac{\partial \widetilde{\mathbf{U}}}{\partial \mathbf{A}(i,j)} = \widetilde{\mathcal{X}}^{\top} \left(\mathbf{I}_{C} \otimes \mathbf{B} \otimes \mathbf{O}_{ij} \right)^{\top} \left(\mathbf{U}_{3} \otimes \mathbf{U}_{2} \otimes \mathbf{U}_{1} \right)^{\top}$$

where $\mathbf{O}_{ij} \in \mathbb{R}^{h \times H}$ is a binary matrix with all entries being zero except the $(i, j)^{\text{th}}$ entry being one. The derivative of ℓ w.r.t. $\widetilde{\mathbf{U}}$ has an explicit form (Yu et al., 2018):

$$\frac{\partial \ell}{\partial \widetilde{\mathbf{U}}} = \widetilde{\mathbf{U}} \left(\mathbf{\Sigma}^{-1} + \mathbf{\Sigma}^{-1} \widetilde{\mathbf{U}}^{\top} \mathbf{y} \mathbf{D}_{\sigma}^{-1} \mathbf{y}^{\top} \widetilde{\mathbf{U}} \mathbf{\Sigma}^{-1} \right) - \mathbf{y} \mathbf{D}_{\sigma}^{-1} \mathbf{y}^{\top} \widetilde{\mathbf{U}} \mathbf{\Sigma}^{-1}$$

where $\Sigma = D_{\sigma} + \widetilde{U}^{\top}\widetilde{U}$. The derivative of \widetilde{U} w.r.t. A, B, U_1, U_2, U_3 can be readily derived by simply replacing each matrix parameter with a sparse binary matrix such as O_{ij} stated above. For example for U_2 , one has:

$$\frac{\partial \widetilde{\mathbf{U}}}{\partial \mathbf{U}_{2}(i,j)} = \widetilde{\mathcal{X}}^{\top} \left(\mathbf{I}_{C} \otimes \mathbf{B} \otimes \mathbf{A} \right)^{\top} \left(\mathbf{U}_{3} \otimes \mathbf{O}_{ij} \otimes \mathbf{U}_{1} \right)^{\top}$$

where $\mathbf{O}_{ij} \in \mathbb{R}^{r_2 \times w}$ and is sparse except the $(i,j)^{\text{th}}$ element being one. The gradients for $\mathbf{U}_1, \mathbf{U}_2, \mathbf{U}_3$ can be used for parameter update, and for \mathbf{A} and \mathbf{B} , we consider updating them via proximal gradient descent. For \mathbf{A} at the $(i+1)^{\text{th}}$ iteration, for example, one applies the gradient descent first to get an estimator proposal: $\widehat{\mathbf{A}}^{(i+\frac{1}{2})} = \widehat{\mathbf{A}}^{(i)} - \eta_i \frac{\partial \ell}{\partial \mathbf{A}}$, and then solves the following optimization problem, which is commonly known as the proximal step:

$$\mathbf{prox}_{\mathsf{TV}}(\widehat{\mathbf{A}}^{(k+\frac{1}{2})}) = \underset{\mathbf{A} \in \mathbb{R}^{h \times H}}{\min} \left\{ \frac{1}{2\eta_i} \left\| \mathbf{A} - \widehat{\mathbf{A}}^{(i+\frac{1}{2})} \right\|_{\mathsf{F}}^2 + \lambda \sum_{s=1}^h \sum_{t=1}^w \| \mathbf{W}_{s,t} \|_{\mathsf{TV}} \right\}$$
(22)

where η_i is the learning rate of the $(i+1)^{\text{th}}$ step.

Solving the proximal problem in (22) can be broken down into multiple parallel 1-D *fused lasso signal approximation* problem. According to Proposition 1 of Friedman et al. (2007), solving (22) can be further broken down into first solving *h* total variation de-noising problem (Rudin et al., 1992):

$$\widetilde{\mathbf{A}}(s,:) \leftarrow \underset{\boldsymbol{\alpha} \in \mathbb{R}^H}{\operatorname{arg \, min}} \, \frac{1}{2\eta_i} \left\| \boldsymbol{\alpha} - \widehat{\mathbf{A}}^{(i+\frac{1}{2})}(s,:) \right\|_{\mathrm{F}}^2 + \lambda \cdot \|\mathbf{B}\|_1 \cdot \sum_{j=2}^H |\boldsymbol{\alpha}(j+1) - \boldsymbol{\alpha}(j)|, \quad s = 1, 2, \dots, h$$
 (23)

Then one can apply a soft-thresholding operator $S_{\lambda\|\nabla_x \mathbf{B}\|_1}(\cdot)$, element-wisely, to $\widetilde{\mathbf{A}} \coloneqq \mathbf{prox}_{\mathrm{TV}}(\widehat{\mathbf{A}}^{(k+\frac{1}{2})})$ to obtain the solution for (22). The problem in (23) can be efficiently solved via the python implementation in $\mathtt{prox}_{\mathrm{TV}}$ based on a fast Newton's method (Jiménez & Sra, 2011; Barbero & Sra, 2018). Similar technique can be applied to update \mathbf{B} , and therefore the final optimization algorithm consists of both a gradient descent step and a fused-lasso proximal step. A more general theoretical discussion on the total variation penalty over 1-D signals can be found in Tibshirani (2014).

The gradient of $\ell(\cdot)$ w.r.t. σ^2 can be easily derived as follows:

$$\frac{\partial \ell}{\partial \sigma^2} = \mathbf{tr} \left[\left(\frac{\partial \ell}{\partial (\mathbf{K} + \mathbf{D}_{\sigma})^{-1}} \right)^{\mathsf{T}} \left(\frac{\partial (\mathbf{K} + \mathbf{D}_{\sigma})^{-1}}{\partial \sigma^2} \right) \right] = \mathbf{tr} \left[\frac{1}{2} \left(\mathbf{K} + \mathbf{D}_{\sigma} \right)^{-1} - \frac{1}{2} \left(\mathbf{K} + \mathbf{D}_{\sigma} \right)^{-2} \mathbf{y} \mathbf{y}^{\mathsf{T}} \right]$$
(24)

Predictions on the unseen testing data with covariates X_* , given the training data (X, y), can be easily derived using the predictive distribution $(y_*|X_*, X, y) \sim \mathcal{N}(\mu_*, \Sigma_*)$:

$$\begin{aligned} \boldsymbol{\mu}_{*} &= \widehat{\mathcal{K}}\left(\mathbf{X}_{*}, \mathbf{X}\right) \left(\widehat{\mathcal{K}}\left(\mathbf{X}, \mathbf{X}\right) + \mathbf{D}_{\widehat{\sigma}}\right)^{-1} \mathbf{y} \\ \boldsymbol{\Sigma}_{*} &= \widehat{\mathcal{K}}\left(\mathbf{X}_{*}, \mathbf{X}_{*}\right) + \mathbf{D}_{\widehat{\sigma}} - \widehat{\mathcal{K}}\left(\mathbf{X}_{*}, \mathbf{X}\right) \left(\widehat{\mathcal{K}}\left(\mathbf{X}, \mathbf{X}\right) + \mathbf{D}_{\widehat{\sigma}}\right)^{-1} \widehat{\mathcal{K}}\left(\mathbf{X}_{*}, \mathbf{X}\right)^{\top} \end{aligned}$$

where $\widehat{\mathcal{K}}(\cdot,\cdot)$ is the kernel function in (7) but evaluated at the estimated model parameters, and $\widehat{\mathcal{K}}(\mathbf{X}_*,\mathbf{X})$ simply denotes

the covariances between the unseen data X_* and the training data X, and the other notations follow.

C. Proof of Theorem 2.2

The proof of Theorem 2.2 is largely based on the convergence results of proximal gradient descent but with additional consideration on the alternating descent scheme. We denote an arbitrary collection of model parameters as $\theta \coloneqq (A, B, U_{1:3}, \sigma)$. Since we update the four blocks of parameters cyclically in Algorithm 1, within each iteration, we further denote the intermediate parameter updates as $\widehat{\theta}^{(k)} \xrightarrow{\text{update } A} \widehat{\theta}^{(k+\frac{1}{4})} \xrightarrow{\text{update } B} \widehat{\theta}^{(k+\frac{1}{2})} \xrightarrow{\text{update } U_{1:3}} \widehat{\theta}^{(k+\frac{3}{4})} \xrightarrow{\text{update } \sigma} \widehat{\theta}^{(k+1)}$. In the remainder of the proof, we will use U to denote $U_{1:3}$ for notational simplicity.

In order to show the upper bound of the difference of the loss function after K iterations with its global minimum $L(\theta^*)$, we first show Lemma C.1:

Lemma C.1. Given the alternating proximal gradient descent algorithm in Algorithm 1 and the assumptions made in Theorem 2.2, one can bound $L(\widehat{\boldsymbol{\theta}}^{(k+\frac{v}{4})}), v \in \{1, 2, 3, 4\}$ as:

$$L(\widehat{\boldsymbol{\theta}}^{(k+\frac{1}{4})}) \le L(\mathbf{A}^*, \widehat{\mathbf{B}}^{(k)}, \widehat{\mathbf{U}}^{(k)}, \widehat{\boldsymbol{\sigma}}^{(k)}) + \frac{1}{2\alpha} \left\{ \left\| \widehat{\mathbf{A}}^{(k)} - \mathbf{A}^* \right\|^2 - \left\| \widehat{\mathbf{A}}^{(k+1)} - \mathbf{A}^* \right\|^2 \right\}$$
(25)

$$L(\widehat{\boldsymbol{\theta}}^{(k+\frac{1}{2})}) \le L(\widehat{\mathbf{A}}^{(k+1)}, \mathbf{B}^{(*)}, \widehat{\mathbf{U}}^{(k)}, \widehat{\boldsymbol{\sigma}}^{(k)}) + \frac{1}{2\alpha} \left\{ \left\| \widehat{\mathbf{B}}^{(k)} - \mathbf{B}^* \right\|^2 - \left\| \widehat{\mathbf{B}}^{(k+1)} - \mathbf{B}^* \right\|^2 \right\}$$
(26)

$$L(\widehat{\boldsymbol{\theta}}^{(k+\frac{3}{4})}) \le L(\widehat{\mathbf{A}}^{(k+1)}, \widehat{\mathbf{B}}^{(k+1)}, \mathbf{U}^*, \widehat{\boldsymbol{\sigma}}^{(k)}) + \frac{1}{2\alpha} \left\{ \left\| \widehat{\mathbf{U}}^{(k)} - \mathbf{U}^* \right\|^2 - \left\| \widehat{\mathbf{U}}^{(k+1)} - \mathbf{U}^* \right\|^2 \right\}$$
(27)

$$L(\widehat{\boldsymbol{\theta}}^{(k+1)}) \le L(\widehat{\mathbf{A}}^{(k+1)}, \widehat{\mathbf{B}}^{(k+1)}, \widehat{\mathbf{U}}^{(k+1)}, \sigma^*) + \frac{1}{2\alpha} \left\{ \left\| \widehat{\sigma}^{(k)} - \sigma^* \right\|^2 - \left\| \widehat{\sigma}^{(k+1)} - \sigma^* \right\|^2 \right\}, \tag{28}$$

where $\|\cdot\|$ is the matrix Frobenious norm, α is a constant learning rate with $\alpha \leq 1/\max\{M_{\mathbf{A}}, M_{\mathbf{B}}, M_{\mathbf{U}}, M_{\sigma}\}$ and $M_{\mathbf{A}}, M_{\mathbf{B}}, M_{\mathbf{U}}, M_{\sigma}$ are the Lipschitz constant for $\mathbf{A}, \mathbf{B}, \mathbf{U}, \sigma$ for the gradient of $\ell(\cdot)$, i.e. the negative log-likelihood defined in (10).

Proof. It suffices to prove (25), and the rest of the bounds follow the same technique.

First, given that the gradient of $\ell(\cdot)$ w.r.t. to \mathbf{A} is Lipschitz continuous with constant $M_{\mathbf{A}}$, i.e. $\|\nabla_{\mathbf{A}}\ell(\mathbf{A}_1) - \nabla_{\mathbf{A}}\ell(\mathbf{A}_2)\| \le M_{\mathbf{A}}\|\mathbf{A}_1 - \mathbf{A}_2\|, \forall \mathbf{A}_1, \mathbf{A}_2$. Since the other parameters also share the same property but have different Lipschitz constant $M_{\mathbf{B}}, M_{\mathbf{U}}, M_{\sigma}$, we use $M := \max\{M_{\mathbf{A}}, M_{\mathbf{B}}, M_{\mathbf{U}}, M_{\sigma}\}$ as the Lipschitz constant for all parameters. Given the Lipschitz continuity of the derivative, one has:

$$\ell(\widehat{\boldsymbol{\theta}}^{(k+\frac{1}{4})}) \le \ell(\widehat{\boldsymbol{\theta}}^{(k)}) + \left\langle \nabla_{\mathbf{A}} \ell(\widehat{\boldsymbol{\theta}}^{(k)}), \widehat{\mathbf{A}}^{(k+1)} - \widehat{\mathbf{A}}^{(k)} \right\rangle + \frac{M}{2} \|\widehat{\mathbf{A}}^{(k+1)} - \widehat{\mathbf{A}}^{(k)}\|^2$$
(29)

which is a direct result from the following inequality for any function $\ell(\cdot)$ with M-Lipschitz continuous derivative:

$$\ell(y) \le \ell(x) + \langle \nabla_x \ell(x), y - x \rangle + \frac{M}{2} ||y - x||^2$$

Additionally, since $\ell(\cdot)$ is assumed as block-wise convex, one has a natural upper bound of $\ell(\widehat{\boldsymbol{\theta}}^{(k)})$ based on convexity:

$$\ell(\widehat{\boldsymbol{\theta}}^{(k)}) \le \ell(\mathbf{A}^*, \widehat{\mathbf{B}}^{(k)}, \widehat{\mathbf{U}}^{(k)}, \widehat{\boldsymbol{\sigma}}^{(k)}) - \left\langle \nabla_{\mathbf{A}} \ell(\widehat{\boldsymbol{\theta}}^{(k)}), \mathbf{A}^* - \widehat{\mathbf{A}}^{(k)} \right\rangle$$
(30)

Combining (29) and (30), one obtains:

$$\ell(\widehat{\boldsymbol{\theta}}^{(k+\frac{1}{4})}) \leq \ell(\mathbf{A}^*, \widehat{\mathbf{B}}^{(k)}, \widehat{\mathbf{U}}^{(k)}, \widehat{\boldsymbol{\sigma}}^{(k)}) + \left\langle \nabla_{\mathbf{A}} \ell(\widehat{\boldsymbol{\theta}}^{(k)}), \widehat{\mathbf{A}}^{(k+1)} - \mathbf{A}^* \right\rangle + \frac{M}{2} \|\widehat{\mathbf{A}}^{(k+1)} - \widehat{\mathbf{A}}^{(k)}\|^2 \tag{31}$$

Also, since $\widehat{\mathbf{A}}^{(k+1)}$ is obtained via a proximal step:

$$\widehat{\mathbf{A}}^{(k+1)} = \operatorname*{arg\,min}_{\mathbf{A}} \frac{1}{2\alpha} \left\| \mathbf{A} - \left(\widehat{\mathbf{A}}^{(k)} - \alpha \nabla_{\mathbf{A}} \ell(\widehat{\boldsymbol{\theta}}^{(k)}) \right) \right\|^2 + \lambda \mathrm{R} \left(\mathbf{A}, \widehat{\mathbf{B}}^{(k)} \right)$$

 $\widehat{\mathbf{A}}^{(k+1)}$ should satisfy the following subgradient condition:

$$G_{\alpha}(\widehat{\boldsymbol{\theta}}^{(k)}) - \nabla_{\mathbf{A}} \ell(\widehat{\boldsymbol{\theta}}^{(k)}) \in \lambda \cdot \partial_{A} R\left(\widehat{\mathbf{A}}^{(k+1)}, \widehat{\mathbf{B}}^{(k)}\right)$$
(32)

where $G_{\alpha}(\widehat{\boldsymbol{\theta}}^{(k)}) \coloneqq -\frac{1}{\alpha} \left(\widehat{\mathbf{A}}^{(k+1)} - \widehat{\mathbf{A}}^{(k)}\right)$ is the proximal gradient. Using the definition of subgradient, one can achieve a trivial inequality as follows:

$$\lambda R\left(\widehat{\mathbf{A}}^{(k+1)}, \widehat{\mathbf{B}}^{(k)}\right) + \left\langle G_{\alpha}(\widehat{\boldsymbol{\theta}}^{(k)}) - \nabla_{\mathbf{A}} \ell(\widehat{\boldsymbol{\theta}}^{(k)}), \mathbf{A}^* - \widehat{\mathbf{A}}^{(k+1)} \right\rangle \le \lambda R\left(\mathbf{A}^*, \widehat{\mathbf{B}}^{(k)}\right)$$
(33)

Combining (31) and (33), we have:

$$\begin{split} L(\widehat{\boldsymbol{\theta}}^{(k+\frac{1}{4})}) &= \ell(\widehat{\boldsymbol{\theta}}^{(k+\frac{1}{4})}) + \lambda R\left(\widehat{\mathbf{A}}^{(k+1)}, \widehat{\mathbf{B}}^{(k)}\right) \\ &\leq L(\mathbf{A}^*, \widehat{\mathbf{B}}^{(k)}, \widehat{\mathbf{U}}^{(k)}, \widehat{\boldsymbol{\sigma}}^{(k)}) + \left\langle G_{\alpha}(\widehat{\boldsymbol{\theta}}^{(k)}), \widehat{\mathbf{A}}^{(k+1)} - \mathbf{A}^* \right\rangle + \frac{M}{2} \|\widehat{\mathbf{A}}^{(k+1)} - \widehat{\mathbf{A}}^{(k)}\|^2 \\ &\leq L(\mathbf{A}^*, \widehat{\mathbf{B}}^{(k)}, \widehat{\mathbf{U}}^{(k)}, \widehat{\boldsymbol{\sigma}}^{(k)}) + \left\langle G_{\alpha}(\widehat{\boldsymbol{\theta}}^{(k)}), \widehat{\mathbf{A}}^{(k)} - \alpha G_{\alpha}(\widehat{\boldsymbol{\theta}}^{(k)}) - \mathbf{A}^* \right\rangle + \frac{1}{2\alpha} \left\| \alpha G_{\alpha}(\widehat{\boldsymbol{\theta}}^{(k)}) \right\|^2 \\ &= L(\mathbf{A}^*, \widehat{\mathbf{B}}^{(k)}, \widehat{\mathbf{U}}^{(k)}, \widehat{\boldsymbol{\sigma}}^{(k)}) + \frac{1}{2\alpha} \left\{ \left\| \widehat{\mathbf{A}}^{(k)} - \mathbf{A}^* \right\|^2 - \left\| \widehat{\mathbf{A}}^{(k+1)} - \mathbf{A}^* \right\|^2 \right\} \end{split}$$

which completes the proof.

In the classical proximal gradient descent context, where one updates a single parameter iteratively, the bound in Lemma C.1 leads to a convergence rate of the algorithm at $\mathcal{O}(1/K)$, after one adds up all the inequalities from iteration 1 to K. The key difference is that, on the right hand side of the inequality (25), the loss function is evaluated at the global minima of \mathbf{A} and the value of \mathbf{B} , \mathbf{U} , σ at the k^{th} iteration. We need to quantify the difference between $L(\mathbf{A}^*, \widehat{\mathbf{B}}^{(k)}, \widehat{\mathbf{U}}^{(k)}, \widehat{\sigma}^{(k)})$ and $L(\boldsymbol{\theta}^*)$ to reach the final error bound result and this difference is given in Lemma C.2 below.

Lemma C.2. The difference of $L(\mathbf{A}^*, \widehat{\mathbf{B}}^{(k)}, \widehat{\mathbf{U}}^{(k)}, \widehat{\sigma}^{(k)})$ and $L(\boldsymbol{\theta}^*) \coloneqq L(\mathbf{A}^*, \mathbf{B}^*, \mathbf{U}^*, \sigma^*)$ can be fully characterized as:

$$L(\mathbf{A}^*, \widehat{\mathbf{B}}^{(k)}, \widehat{\mathbf{U}}^{(k)}, \widehat{\sigma}^{(k)}) - L(\boldsymbol{\theta}^*) \le \frac{M}{2} \left\{ \|\widehat{\mathbf{B}}^{(k)} - \mathbf{B}^*\| + \|\widehat{\mathbf{U}}^{(k)} - \mathbf{U}^*\| + \|\widehat{\sigma}^{(k)} - \sigma^*\| \right\}^2$$
(34)

+
$$\|\nabla_{\mathbf{B}}\ell(\boldsymbol{\theta}^*)\| \cdot \|\widehat{\mathbf{B}}^{(k)} - \mathbf{B}^*\| + \lambda R\left(\mathbf{A}^*, \widehat{\mathbf{B}}^{(k)} - \mathbf{B}^*\right)$$
 (35)

where (34) is the additional loss incurred by using the iterative value of the other parameters instead of the global optimum (called the ALT-gap), and (35) is the additional loss incurred by using the total variation penalty with the $k^{\rm th}$ iterative value of ${\bf B}$ (called the TV-gap).

Proof. We start the derivation with the following trivial decomposition:

$$L(\mathbf{A}^*, \widehat{\mathbf{B}}^{(k)}, \widehat{\mathbf{U}}^{(k)}, \widehat{\sigma}^{(k)}) - L(\boldsymbol{\theta}^*) = L(\mathbf{A}^*, \widehat{\mathbf{B}}^{(k)}, \widehat{\mathbf{U}}^{(k)}, \widehat{\sigma}^{(k)}) - L(\mathbf{A}^*, \mathbf{B}^*, \widehat{\mathbf{U}}^{(k)}, \widehat{\sigma}^{(k)})$$
(36)

$$+L(\mathbf{A}^*, \mathbf{B}^*, \widehat{\mathbf{U}}^{(k)}, \widehat{\sigma}^{(k)}) - L(\mathbf{A}^*, \mathbf{B}^*, \mathbf{U}^*, \widehat{\sigma}^{(k)})$$
(37)

$$+L(\mathbf{A}^*, \mathbf{B}^*, \mathbf{U}^*, \widehat{\sigma}^{(k)}) - L(\boldsymbol{\theta}^*)$$
(38)

We need to bound (36), (37) and (38) separately. For (36), we have:

$$L(\mathbf{A}^{*}, \widehat{\mathbf{B}}^{(k)}, \widehat{\mathbf{U}}^{(k)}, \widehat{\sigma}^{(k)}) - L(\mathbf{A}^{*}, \mathbf{B}^{*}, \widehat{\mathbf{U}}^{(k)}, \widehat{\sigma}^{(k)})$$

$$\leq \ell(\mathbf{A}^{*}, \widehat{\mathbf{B}}^{(k)}, \widehat{\mathbf{U}}^{(k)}, \widehat{\sigma}^{(k)}) - \ell(\mathbf{A}^{*}, \mathbf{B}^{*}, \widehat{\mathbf{U}}^{(k)}, \widehat{\sigma}^{(k)}) + \lambda \mathbf{R} \left(\mathbf{A}^{*}, \widehat{\mathbf{B}}^{(k)} - \mathbf{B}^{*}\right)$$

$$\leq \left\langle \nabla_{\mathbf{B}} \ell(\mathbf{A}^{*}, \mathbf{B}^{*}, \widehat{\mathbf{U}}^{(k)}, \widehat{\sigma}^{(k)}), \widehat{\mathbf{B}}^{(k)} - \mathbf{B}^{*} \right\rangle + \frac{M}{2} \|\widehat{\mathbf{B}}^{(k)} - \mathbf{B}^{*}\|^{2} + \lambda \mathbf{R} \left(\mathbf{A}^{*}, \widehat{\mathbf{B}}^{(k)} - \mathbf{B}^{*}\right)$$

$$\leq \|\widehat{\mathbf{B}}^{(k)} - \mathbf{B}^{*}\| \cdot \left(\|\nabla_{\mathbf{B}} \ell(\boldsymbol{\theta}^{*})\| + M \|\widehat{\mathbf{U}}^{(k)} - \mathbf{U}^{*}\| + M \|\widehat{\boldsymbol{\sigma}}^{(k)} - \sigma^{*}\|\right) + \frac{M}{2} \|\widehat{\mathbf{B}}^{(k)} - \mathbf{B}^{*}\|^{2} + \lambda \mathbf{R} \left(\mathbf{A}^{*}, \widehat{\mathbf{B}}^{(k)} - \mathbf{B}^{*}\right)$$

$$(39)$$

where the last line follows from the Cauchy-Schwartz inequality followed by the Lipschitz continuity of the gradient and the triangle inequality of the Frobenius norm.

As for (37), similar technique follows and lead to:

$$\begin{split} L(\mathbf{A}^*, \mathbf{B}^*, \widehat{\mathbf{U}}^{(k)}, \widehat{\boldsymbol{\sigma}}^{(k)}) - L(\mathbf{A}^*, \mathbf{B}^*, \mathbf{U}^*, \widehat{\boldsymbol{\sigma}}^{(k)}) &= \ell(\mathbf{A}^*, \mathbf{B}^*, \widehat{\mathbf{U}}^{(k)}, \widehat{\boldsymbol{\sigma}}^{(k)}) - \ell(\mathbf{A}^*, \mathbf{B}^*, \mathbf{U}^*, \widehat{\boldsymbol{\sigma}}^{(k)}) \\ &\leq \left\langle \nabla_{\mathbf{U}} \ell(\mathbf{A}^*, \mathbf{B}^*, \mathbf{U}^*, \widehat{\boldsymbol{\sigma}}^{(k)}), \widehat{\mathbf{U}}^{(k)} - \mathbf{U}^* \right\rangle + \frac{M}{2} \|\widehat{\mathbf{U}}^{(k)} - \mathbf{U}^*\|^2 \\ &\leq M \|\widehat{\mathbf{U}}^{(k)} - \mathbf{U}^*\| \cdot \|\widehat{\boldsymbol{\sigma}}^{(k)} - \boldsymbol{\sigma}^*\| + \frac{M}{2} \|\widehat{\mathbf{U}}^{(k)} - \mathbf{U}^*\|^2 \end{split}$$

where the last line uses the fact that at the global optimum, we have $\nabla_{\mathbf{U}}\ell(\boldsymbol{\theta}^*) = 0$.

Similarly for (38), one has:

$$L(\mathbf{A}^*, \mathbf{B}^*, \mathbf{U}^*, \widehat{\sigma}^{(k)}) - L(\boldsymbol{\theta}^*) \le \left\langle \nabla_{\sigma} \ell(\boldsymbol{\theta}^*), \widehat{\sigma}^{(k)} - \sigma^* \right\rangle + \frac{M}{2} \|\widehat{\sigma}^{(k)} - \sigma^*\|^2$$

Combining the three individual upper bounds together yields the result and thereby completes the proof.

Similar results in Lemma C.2 can be easily derived for \mathbf{B}^* , \mathbf{U}^* , σ^* . With these theoretical results, we are now ready to prove Theorem 2.2:

Proof. Combining the results in Lemma C.1 and C.2, we have the following upper bound for $L(\widehat{\boldsymbol{\theta}}^{(k+\frac{1}{4})}) - L(\boldsymbol{\theta}^*)$:

$$L(\widehat{\boldsymbol{\theta}}^{(k+\frac{1}{4})}) - L(\boldsymbol{\theta}^*) \leq \frac{M}{2} \left\{ \|\widehat{\mathbf{B}}^{(k)} - \mathbf{B}^*\| + \|\widehat{\mathbf{U}}^{(k)} - \mathbf{U}^*\| + \|\widehat{\sigma}^{(k)} - \sigma^*\| \right\}^2$$

$$+ \|\nabla_{\mathbf{B}}\ell(\boldsymbol{\theta}^*)\| \cdot \|\widehat{\mathbf{B}}^{(k)} - \mathbf{B}^*\| + \lambda \mathbf{R} \left(\mathbf{A}^*, \widehat{\mathbf{B}}^{(k)} - \mathbf{B}^* \right)$$

$$+ \frac{1}{2\alpha} \left\{ \left\| \widehat{\mathbf{A}}^{(k)} - \mathbf{A}^* \right\|^2 - \left\| \widehat{\mathbf{A}}^{(k+1)} - \mathbf{A}^* \right\|^2 \right\}$$

If $\mathbf{B}, \mathbf{U}, \sigma$ is fixed at $\mathbf{B}^*, \mathbf{U}^*, \sigma^*$ and one only updates \mathbf{A} via the proximal gradient descent, the error bound here vanishes to its last term only and can be further reduced if one adds up the inequality from iteration 1 to K, leading to the classical proximal gradient descent convergence rate result. Similar bounds for $L(\widehat{\boldsymbol{\theta}}^{(k+\frac{1}{2})}), L(\widehat{\boldsymbol{\theta}}^{(k+\frac{3}{4})}), L(\widehat{\boldsymbol{\theta}}^{(k+1)})$ can be derived

and we aggregate the results together as follows:

$$\sum_{k=0}^{K} \sum_{v=1}^{4} \left(L(\widehat{\boldsymbol{\theta}}^{(k+\frac{v}{4})}) - L(\boldsymbol{\theta}^{*}) \right) \leq \frac{1}{2\alpha} \left\{ \|\widehat{\boldsymbol{\theta}}^{(0)} - \boldsymbol{\theta}^{*}\|^{2} - \|\widehat{\boldsymbol{\theta}}^{(K)} - \boldsymbol{\theta}^{*}\|^{2} \right\}
+ \sum_{k=0}^{K} h_{\lambda} (\widehat{\mathbf{A}}^{(k+1)} - \mathbf{A}^{*}, \widehat{\mathbf{B}}^{(k+1)} - \mathbf{B}^{*}, \widehat{\mathbf{B}}^{(k)} - \mathbf{B}^{*})
+ \frac{1}{2\alpha} \sum_{k=0}^{K} \tau \left(\widehat{\boldsymbol{\theta}}^{(k+1)}, \widehat{\boldsymbol{\theta}}^{(k)}, \boldsymbol{\theta}^{*} \right)$$
(40)

where $h_{\lambda}(\widehat{\mathbf{A}}^{(k+1)} - \mathbf{A}^*, \widehat{\mathbf{B}}^{(k+1)} - \mathbf{B}^*, \widehat{\mathbf{B}}^{(k)} - \mathbf{B}^*)$ is the extra gap from the optimal loss created by the existence of the total variation penalty in the loss function (thus we call it the TV-gap) and is defined as:

$$h_{\lambda}(\widehat{\mathbf{A}}^{(k+1)} - \mathbf{A}^{*}, \widehat{\mathbf{B}}^{(k+1)} - \mathbf{B}^{*}, \widehat{\mathbf{B}}^{(k)} - \mathbf{B}^{*}) := 3\|\nabla_{\mathbf{A}}\ell(\boldsymbol{\theta}^{*})\| \cdot \|\widehat{\mathbf{A}}^{(k+1)} - \mathbf{A}^{*}\| + 2\|\nabla_{\mathbf{B}}\ell(\boldsymbol{\theta}^{*})\| \cdot \|\widehat{\mathbf{B}}^{(k+1)} - \mathbf{B}^{*}\| + \|\nabla_{\mathbf{B}}\ell(\boldsymbol{\theta}^{*})\| \cdot \|\widehat{\mathbf{B}}^{(k)} - \mathbf{B}^{*}\| + \lambda R\left(\mathbf{A}^{*}, \widehat{\mathbf{B}}^{(k)} - \mathbf{B}^{*}\right) + 3\lambda R\left(\widehat{\mathbf{A}}^{(k+1)} - \mathbf{A}^{*}, \mathbf{B}^{*}\right) + 2\lambda R\left(\mathbf{A}^{*}, \widehat{\mathbf{B}}^{(k+1)} - \mathbf{B}^{*}\right) + 2\lambda R\left(\widehat{\mathbf{A}}^{(k+1)} - \mathbf{A}^{*}, \widehat{\mathbf{B}}^{(k+1)} - \mathbf{B}^{*}\right)$$

$$(41)$$

where $R(\mathbf{A}, \mathbf{B})$ is the total variation penalty defined in Lemma 2.1. $\tau(\widehat{\boldsymbol{\theta}}^{(k+1)}, \widehat{\boldsymbol{\theta}}^{(k)}, \boldsymbol{\theta}^*)$ is the extra gap from the optimal loss created by the usage of iterative value of the parameters during the alternating proximal gradient descent (thus we call it the ALT-gap) and is defined as:

$$\tau\left(\widehat{\boldsymbol{\theta}}^{(k+1)}, \widehat{\boldsymbol{\theta}}^{(k)}, \boldsymbol{\theta}^{*}\right) := \left[\|\widehat{\mathbf{B}}^{(k)} - \mathbf{B}^{*}\| + \|\widehat{\mathbf{U}}^{(k)} - \mathbf{U}^{*}\| + \|\widehat{\boldsymbol{\sigma}}^{(k)} - \boldsymbol{\sigma}^{*}\|\right]^{2} \\
+ \left[\|\widehat{\mathbf{A}}^{(k+1)} - \mathbf{A}^{*}\| + \|\widehat{\mathbf{U}}^{(k)} - \mathbf{U}^{*}\| + \|\widehat{\boldsymbol{\sigma}}^{(k)} - \boldsymbol{\sigma}^{*}\|\right]^{2} \\
+ \left[\|\widehat{\mathbf{A}}^{(k+1)} - \mathbf{A}^{*}\| + \|\widehat{\mathbf{B}}^{(k+1)} - \mathbf{B}^{*}\| + \|\widehat{\boldsymbol{\sigma}}^{(k)} - \boldsymbol{\sigma}^{*}\|\right]^{2} \\
+ \left[\|\widehat{\mathbf{A}}^{(k+1)} - \mathbf{A}^{*}\| + \|\widehat{\mathbf{B}}^{(k+1)} - \mathbf{B}^{*}\| + \|\widehat{\mathbf{U}}^{(k+1)} - \mathbf{U}^{*}\|\right]^{2} \tag{42}$$

The final result can be derived from (40) by lower bounding the left hand side:

$$\sum_{k=0}^K \sum_{v=1}^4 \left(L(\widehat{\boldsymbol{\theta}}^{(k+\frac{v}{4})}) - L(\boldsymbol{\theta}^*) \right) \geq 4(K+1) \left(L(\widehat{\boldsymbol{\theta}}^{(K+1)}) - L(\boldsymbol{\theta}^*) \right)$$

which is evident given that each step is a descent step.

Although we cannot fully verify the assumptions made, we plot the history of the loss function and the relative change of the model parameters for our real data application in Figure 5. Empirically, our model demonstrates a convergence rate at $\mathcal{O}(1/K)$ (see the red curve fitted based on a polynomial model with function form $f(k) = a + \frac{b}{c+k}$).

D. Details of the Simulation Study

In this section we provide the details on generating the simulation data. Given the three types of data in Figure 2, we use two sparse and banded tensor contracting factors $(\mathbf{A}^*, \mathbf{B}^*)$ (see the top of Figure 6a) to contract each channel to a 3×3 tensor (see Figure 1a about the contraction operation). $(\mathbf{A}^*, \mathbf{B}^*)$ in Figure 6a essentially do a 5×5 block averaging for the four corners, four sides and the middle block of each channel data. So one can expect the Type 1 & 3 data in Figure 2 to have its *signal* in the four corners of the contracted tensor, and Type 2 has its *signal* in the middle block (see Figure 2 for an illustration). Given the contracted tensor, we use the multi-linear kernel, specified in Figure 6a (bottom) to generate the

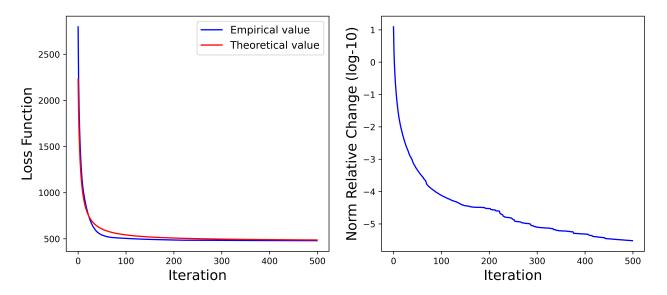


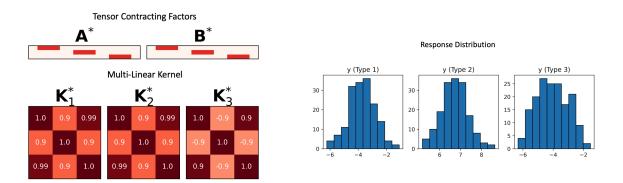
Figure 5: (Left) Loss history of the solar flare intensity regression task with Tensor-GPST ($\lambda=1.0$); A curve at the order of $\mathcal{O}(1/K)$ is fitted to the loss history and empirically, the algorithm converges at the rate of $\mathcal{O}(1/K)$ to a local minimum. (Right) History of the Frobenius norm of the relative change of model parameters in log-10 scale, which suggests that the parameters converge to a stationary point, and thus the ALT-gap and the TV-gap will converge to a constant.

response variables via:

$$\mathbf{y} \sim \mathcal{N}\left(\mathbf{0}_N, \mathbf{K}^* + \sigma^2 \mathbf{I}_N\right)$$

where
$$\mathbf{K}^*(i,j) = \mathbf{vec}(\mathcal{X}_i)^{\top} \left[\mathbf{K}_3^* \otimes \left(\mathbf{B}^{\top} \mathbf{K}_2^* \mathbf{B} \right) \otimes \left(\mathbf{A}^{\top} \mathbf{K}_1^* \mathbf{A} \right) \right] \mathbf{vec}(\mathcal{X}_j)$$
 and $\sigma = 0.5$.

One can notice from the kernel \mathbf{K}_3^* in Figure 6a that channel 1 & 3 are positively correlated and channel 2 is negatively correlated with both channel 1 & 3, and this is reflected in Figure 6b, where we plot the distribution of the simulated sample of size N=500, by the type of data. The tensor regression problem is to use the original $25 \times 25 \times 3$ tensors \mathcal{X}_i to forecast the regression label y_i .



(a) True tensor contracting factors (top) and true multi-linear ker-(b) Distribution of the response variable y by type (see type defining (bottom). tion in Figure 2. Total sample size N=500.

Figure 6: Ground Truth of the Simulated data. (a) The true tensor contracting factors $(\mathbf{A}^*, \mathbf{B}^*)$ (top), where each has a banded structure with the 5 consecutive pixels filled with 0.2 on each row. The bottom shows the multi-linear kernel $\mathbf{K}_1^*, \mathbf{K}_2^*, \mathbf{K}_3^*$. (b) The resulting response distribution of each type of data. One can see how type 1 & 3 has similar distribution, thanks to their high channel correlation in \mathbf{K}_3^* .

E. Details of the AIA-HMI Solar Flare Imaging Dataset

In this appendix, we provide some astrophysical backgrounds and additional details on data preprocessing about the AIA-HMI solar flare imaging datasets.

There are over 12,000 solar flares recorded by the Geostationary Operational Environmental Satellite (GOES) from May, 2010 to June, 2017, with intensity at least at the A-class flare level (peak X-ray brightness $< 10^{-7} \text{W/m}^2$). Among these flares, 4,409 are B-class flares ($10^{-7} \sim 10^{-6} \text{W/m}^2$), 710 are M-class flares ($10^{-5} \sim 10^{-4} \text{W/m}^2$) and 50 are X-class flares ($> 10^{-4} \text{W/m}^2$). We combine the M-class and X-class flares in a single class, we name the class as M/X-class flares. Each flare is associated with a solar active region, which is a localized, transient volume of the solar atmosphere characterized by complex magnetic fields. We collect the AIA and HMI imaging data for each of the M/X-class flare during this period, and collect the B-class flares happened within the same active regions to construct our own database. Given the data availability, we end up with a database of 1,264 B-class flares and 728 M/X-class flares.

The AIA imaging data has 8 channels, distinguished by the wavelength band of the Extreme Ultraviolet (EUV) and Ultraviolet (UV) spectrum used to image the Sun 3 . The AIA channels are named under their respective spectral band: AIA-94Å,AIA-131Å, AIA-171Å, AIA-193Å, AIA-211Å, AIA-304Å, AIA-335Å and AIA-1600Å. The HMI imaging data captures the r, θ, ϕ -component of the solar magnetic field, and in our database, we keep the HMI B_r channel, which have demonstrated contains flare-predictive signals (Sun et al., 2022). Finally, we derive the polarity inversion line (PIL) (Schrijver, 2007) from the B_r , which highlights a sub-region with the strongest flare discriminating signals (Wang et al., 2020; Sun et al., 2021) and $B_r \approx 0$. In Figure 7, we plot one example of the 10 channels for an M-class flare.

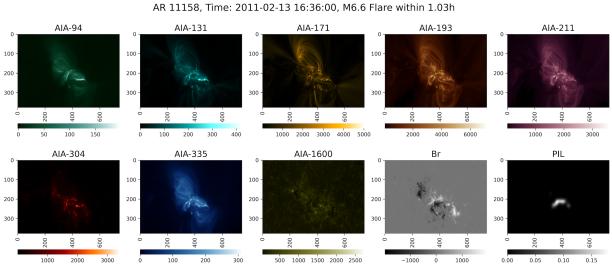


Figure 7: M-class Flare Example for Active Region (AR) No.11158, recorded at 16:36:00 (UT) of Feb 13, 2021. The flare intensity is $6.6 \times 10^{-5} \text{W/m}^2$ and peaked at 17:38:00 (UT) of the same day. Tensor data size is $377 \times 744 \times 10$. Channel name labeled on top of each panel where we omit the Å.

For the particular case in Figure 7, the image size is 377×744 , but different active regions are of different size. Also, different flares have their PIL, as well as the major signals in the other channels, stretching in different directions. To unify the size and orientation of all flares' imaging data, we follow these steps to preprocess our data:

- Pick the pixel in the PIL channel with the largest sum of PIL weights near its 51×51 neighborhood. This helps on picking the "center" of the image. If the PIL only contains zeros, which could happen for some very weak B-class flares, we use the AIA-1600Å in place of the PIL and do the same thing.
- Around the "center", we randomly sample 5,000 pixels, with replacement and the sampling probability is proportional to the PIL (or AIA-1600Å) pixel intensity, and do a Principal Component Analysis (PCA) of each pixel's 2D (x, y)

³See more details at https://sdo.gsfc.nasa.gov/data/channels.php

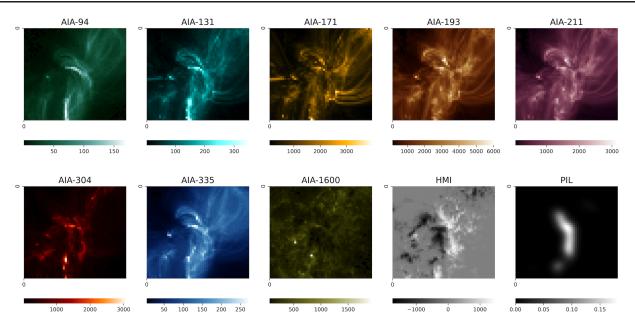


Figure 8: Pre-processed version of the sample in Figure 7. Notice how the PIL channel is now aligned vertically. Tensor size is reduced to $50 \times 50 \times 10$ for all 1,329 flares.

coordinates (coordinates on the pixel grid) and use the first principal component to calculate the orientation of the PIL. This step helps to find the "direction" of the image.

• We rotate each channel with the same angle such that the "direction" of the PIL is vertical. Then, we crop a 201×201 window around the "center" of the image, and do zero-padding where it is needed.

These preprocessing steps create flare data that are roughly comparable, but just as the simulation data pattern in Figure 2, there is still randomness w.r.t. the positioning and direction of the flare predictive signal for each individual flare. We subset our flare list to those whose longitude is within $\pm 60^{\circ}$ from the Sun's central meridian, which removes the low-quality samples with limb distortion. This reduces our sample size from 1,992 flares to 1,329 flares. We further reduces the dimensionality of the 201×201 images to 50×50 , after applying the preprocessing steps above, by bi-linear interpolation to speed up the model computation. The pre-processed version of the sample in Figure 7, with tensor size $50 \times 50 \times 10$, is shown in Figure 8. Notice how the PIL channel is now looking more "vertical" and how each channel is sort of "zoomed-in". The tensor size is now unified across all samples as $50 \times 50 \times 10$.

Before running the model, we normalize the scale of each channel such that each channel has its pixel intensity roughly within the range of [-1, 1], to avoid numerical overflow in the algorithm. We only use the training set scale information to determine the scaling factor in order to avoid information spillover.

For the flare intensity, originally, B-class flare has its intensity within the interval $[10^{-7}, 10^{-6}]$ (unit: W/m²), and M/X-class flare has its intensity within the interval $[10^{-5}, +\infty]$ (unit: W/m²). We transform any intensity y of each flare via:

$$\tilde{y} = \log_{10}(y) + 5.5$$

such that the middle point of the weakest M/X-class flare and the strongest B-class flare is centered at zero.

F. Additional Results on Solar Flare Forecasting

This appendix provides additional results on the solar flare intensity regression. We first visualize the parameter estimates of **GPST**, which is the best-performing model in Table 3, under one random train/test split with $\lambda=1.0$ in this section. Figure 9 provides the kernel estimates (the left three panels) and Figure 10 shows the non-zero feature maps.

The kernel estimators $\hat{\mathbf{K}}_1$ and $\hat{\mathbf{K}}_2$ indicate that feature map $\mathbf{W}_{1,2}$ is of great importance since $\hat{\mathbf{K}}_1(1,1)$ and $\hat{\mathbf{K}}_2(2,2)$

contains the largest element, indicating that the feature extracted by $W_{1,2}$ explains the most variations across all feature maps.

To formally conceptualize the notion of feature map importance as well as channel importance, one can start by decomposing the variations of the regression label y given tensor data $\mathcal{X} \in \mathbb{R}^{H \times W \times C}$ as follows:

$$\operatorname{Var}(y) = \sum_{\substack{1 \leq s_1, s_2 \leq h \\ 1 \leq t_1, t_2 \leq w \\ 1 < c_1, c_2 \leq C}} \underbrace{\mathbf{K}_1(s_1, s_2) \cdot \mathbf{K}_2(t_1, t_2)}_{\text{Feature Map Importance}} \times \underbrace{\mathbf{K}_3(c_1, c_2)}_{\text{Channel Importance}} \times \underbrace{\left\langle \mathbf{W}_{s_1, t_1}, \mathbf{X}^{(c_1)} \right\rangle \cdot \left\langle \mathbf{W}_{s_2, t_2}, \mathbf{X}^{(c_2)} \right\rangle}_{\text{Latent Features Product}} + \underbrace{\sigma^2}_{\text{Capper Channel Importance}}, \tag{43}$$

and this leads to a natural definition of the percentage of explained variation for any pair of channels $c_1, c_2 \in [C]$:

$$\% \text{ Explained Variation} = \frac{\mathbf{K}_3(c_1, c_2)}{\text{Var}(y)} \times \sum_{\substack{1 \le s_1, s_2 \le h \\ 1 \le t_1, t_2 \le w}} \left\langle \mathbf{W}_{s_1, t_1}, \mathbf{X}^{(c_1)} \right\rangle \cdot \left\langle \mathbf{W}_{s_2, t_2}, \mathbf{X}^{(c_2)} \right\rangle \times \mathbf{K}_1(s_1, s_2) \cdot \mathbf{K}_2(t_1, t_2) \quad (44)$$

Similarly, one can define the percentage of explained variation for any pair of feature maps \mathbf{W}_{s_1,t_1} , \mathbf{W}_{s_2,t_2} , with $s_1,s_2 \in [h],t_1,t_2 \in [w]$, as:

$$\% \text{ Explained Variation} = \frac{\mathbf{K}_1(s_1, s_2) \times \mathbf{K}_2(t_1, t_2)}{\text{Var}(y)} \times \sum_{1 \leq c_1, c_2 \leq C} \left\langle \mathbf{W}_{s_1, t_1}, \mathbf{X}^{(c_1)} \right\rangle \cdot \left\langle \mathbf{W}_{s_2, t_2}, \mathbf{X}^{(c_2)} \right\rangle \times \mathbf{K}_3(c_1, c_2) \tag{45}$$

The analysis here is a by-product of the Tensor-GPST model and is similar to the Joint and Individual Variation Explained (JIVE) (Lock et al., 2013) analysis. Both (44) and (45) can be computed empirically by plugging in the parameter estimators of $\hat{\mathbf{K}}_1, \hat{\mathbf{K}}_2, \hat{\mathbf{K}}_3$ and $\hat{\mathbf{W}}_{s,t}, s \in [h], t \in [w]$ and use all training inputs \mathcal{X} for calculation and take an average.

In the last two panels of Figure 9, we show the percentage of explained variation for all 10 AIA-HMI channels based on (44) and the percentage of explained variation for all 9 feature maps based on (45). For channel-wise explained variation, we simply fix $c_1 = c_2$ in (44), and for feature map explained variation, we simply fix $(s_1, t_1) = (s_2, t_2)$. Note that since all channels share the same set of feature maps, the latent features of different channel are not orthogonal, which indicates that the sum of the percentage of explained variation defined in (44) could exceed 100%. The same argument holds for the feature maps' explained variation. But the explained variation still reveals the relative importance of channels and feature maps.

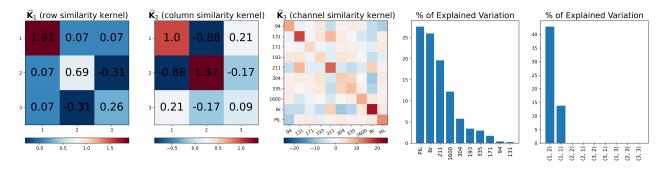


Figure 9: **GPST** (under random train/test split, $\lambda = 1.0$) kernel estimates (panel 1-3), channel-wise % of explained variations (panel 4) and feature map % of explained variations (panel 5). It coincides with the literature (Wang et al., 2020; Sun et al., 2021) that the PIL is the channel with strong flare signals and the AIA imaging data is a good add-on to the HMI channel. The index for feature maps are the 2-tuple (s,t).

The feature maps shown in Figure 10 mainly highlight two patterns:

• All six feature maps show non-zero weights on at least one of the four boundaries. This indicates that the features collected are around the perimeter of the flare eruptive region, which captures the "size" of the flare eruptive area. In

Figure 11 and 12, we show the sample average of all 10 channels for the M/X-class and B-class flares, respectively. One can easily notice the difference between the two classes in terms of the "size" of the bright spots.

• There are some non-zero weights in $W_{1,2}$ and other feature maps near row 20, where features are collected near the top of the brightest PIL region of the M/X flares.

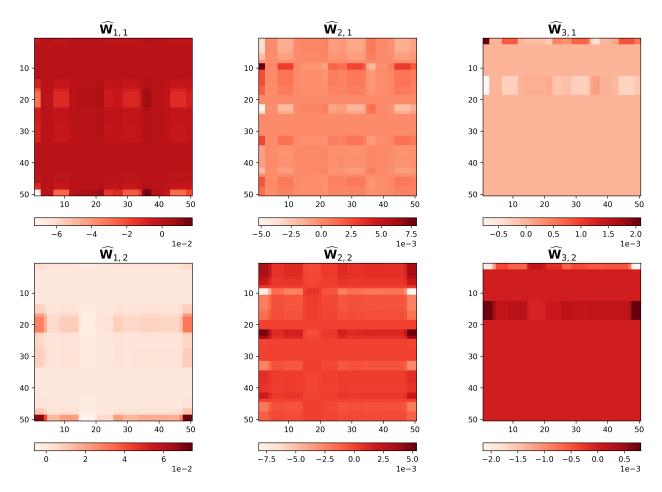


Figure 10: **GPST** (under random train/test split, $\lambda = 1.0$) feature map (the non-zero ones) estimates.

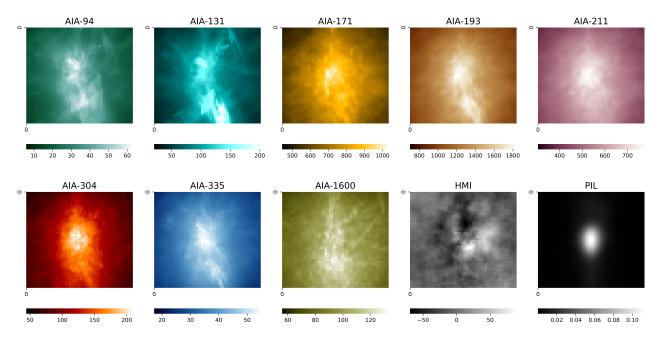


Figure 11: Sample average AIA-HMI map for M-class flare.

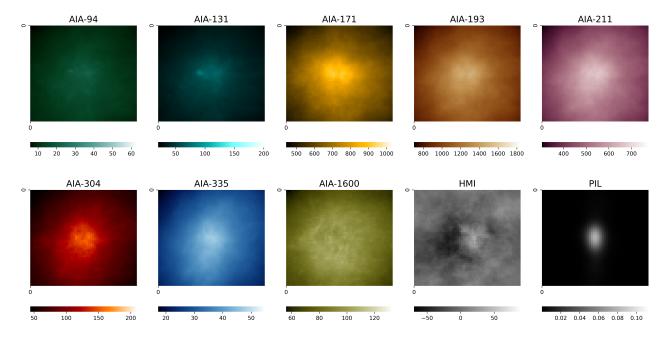


Figure 12: Sample average AIA-HMI map for B-class flare.