



# Fast Estimation of the Kernel Group LASSO

Zoltán Ádám Milacski\*, Son Truong Hy, Balázs Pintér and András Lőrincz

Neural Information Processing Group, Department of Software Technology and Methodology, Faculty of Informatics, Eötvös Loránd University, Hungary



## Abstract

The Kernel Group LASSO is an  $\ell_1/\ell_2$  regularized (structured sparse)  $\ell_2$  reconstruction problem, which performs well at multi-label classification and is defined in a Reproducing Kernel Hilbert Space. Unfortunately, computing the ground truth solution to this task is slow for real-time applications even with state-of-the-art optimization schemes like the Fast Iterative Shrinkage Thresholding Algorithm. We extend the Learned Iterative Shrinkage Thresholding Algorithm – a fast neural network introduced by Gregor and LeCun – to estimate the true result. We test our method in time series classification by training on the 6D Motion Gesture Database while utilizing the Global Alignment time series kernel.

## Kernel Group LASSO

Denote by  $\mathbf{X} \neq \emptyset$  a set where the symmetric, positive semidefinite normalized kernel function  $\mathbf{k}: \mathbf{X} \times \mathbf{X} \rightarrow \mathbb{R}$  is defined:

$$\mathbf{k}(\mathbf{x}, \mathbf{y}) = \frac{\langle \varphi(\mathbf{x}), \varphi(\mathbf{y}) \rangle_{\mathcal{H}}}{\sqrt{\langle \varphi(\mathbf{x}), \varphi(\mathbf{x}) \rangle_{\mathcal{H}}} \sqrt{\langle \varphi(\mathbf{y}), \varphi(\mathbf{y}) \rangle_{\mathcal{H}}}}, \quad \forall (\mathbf{x}, \mathbf{y}) \in \mathbf{X} \times \mathbf{X} \quad (1)$$

for Reconstructing Kernel Hilbert Space  $\mathcal{H}$  and feature mapping  $\varphi: \mathbf{X} \rightarrow \mathcal{H}$ . Let  $\mathbf{x} \in \mathbf{X}$  and  $\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_N] \in \mathbf{X}$  be a signal and a vector system (i.e., a dictionary), respectively. Then for group structure  $\mathcal{G} \subseteq 2^{\{1, \dots, N\}}, \cup_{\mathbf{G} \in \mathcal{G}} \mathbf{G} = \{1, \dots, N\}$ ,  $\lambda > 0$ , consider the **Kernel Group LASSO (KGLASSO)** procedure [1]:

$$\alpha^* = \arg \min_{\alpha \in \mathbb{R}^N} \frac{1}{2} \alpha^T \mathbf{k}(\mathbf{D}, \mathbf{D}) \alpha - \mathbf{k}(\mathbf{D}, \mathbf{x})^T \alpha + \lambda \sum_{\mathbf{G} \in \mathcal{G}} \left( \sqrt{|\mathbf{G}|} \cdot \|\alpha_{\mathbf{G}}\|_2 \right), \quad (2)$$

i.e., we aim to **reconstruct  $\varphi(\mathbf{x})$  with group sparse linear combination of  $\varphi(\mathbf{D})$**  within  $\mathcal{H}$ . This is an  $\ell_1/\ell_2$  regularized quadratic programming problem that can be solved by the Fast Iterative Shrinkage Thresholding Algorithm (FISTA) [2] after precomputing  $\mathbf{k}(\mathbf{D}, \mathbf{D})$  and  $\mathbf{k}(\mathbf{D}, \mathbf{x})$ .

### Advantages:

- Kernels can discover *more subtle similarities* and make the system undercomplete.
- Group structure can *reduce the problem size* by choosing from fewer variables.
- Normalized group activations  $\left( \frac{\|\alpha_{\mathbf{G}}^*\|_2}{\sqrt{|\mathbf{G}|}} \right)_{\mathbf{G} \in \mathcal{G}}$  can be used for *multi-label classification*.

### Limitations:

- FISTA still has significant *time complexity* as it requires several iterations.
- KGLASSO *scales quadratically in dictionary size  $N$* .
- KGLASSO *scales linearly in signal count*.
- Kernel computations can be very slow and often lead to *dense matrices*, which are difficult to store and deal with.

## Learned Iterative Shrinkage Thresholding Algorithm

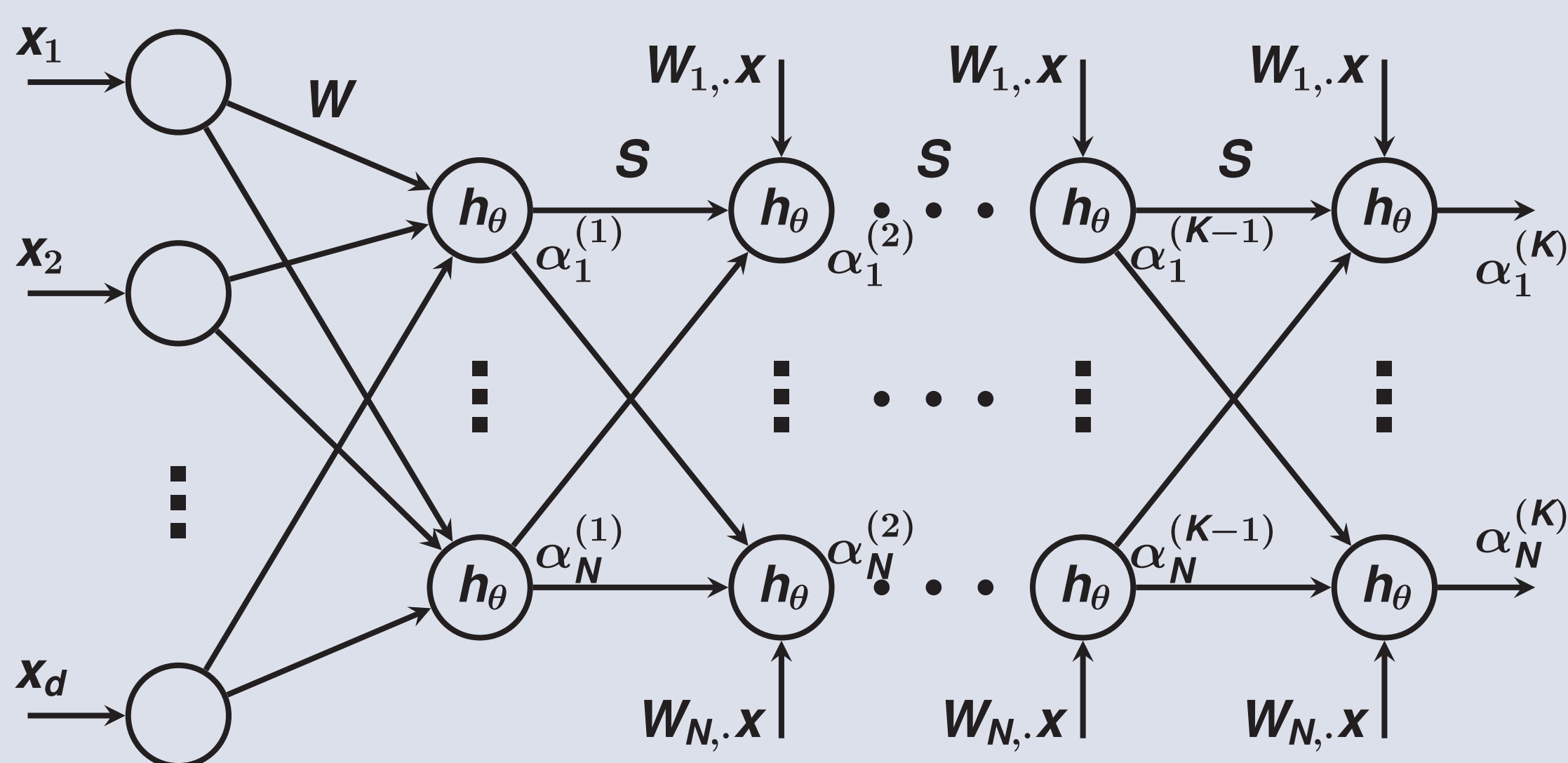
Due to the FISTA limitation above, supervised approximation schemes to it have gained attention recently. The **Learned Iterative Shrinkage Thresholding Algorithm (LISTA)** was proposed in [3]: a **neural network for estimating the sparse code** of the original linear and unstructured LASSO [4], where:

- *adaptive soft-thresholding activation function* was introduced to yield true sparse outputs with respect to tunable thresholds  $\theta \in \mathbb{R}^N$ :

$$\mathbf{h}_{\theta}(\alpha) = \text{sign}(\alpha) \circ (|\alpha| - \theta)_+, \quad (3)$$

- *competition between dictionary elements* was introduced by making the network recurrent with fixed depth  $\mathbf{K} \in \mathbb{N}$ :

$$\alpha^{(0)} = \mathbf{0}, \quad \alpha^{(k)} = \mathbf{h}_{\theta} \left( \mathbf{W} \mathbf{x} + \mathbf{S} \alpha^{(k-1)} \right), \quad k = 1, \dots, \mathbf{K}, \quad (4)$$



i.e., signal  $\mathbf{x}$  is mapped to code space right away with  $\mathbf{W}$  and then  $\mathbf{S}$  rules out some of the active elements,

- *minimization of  $\ell_2$  loss* was carried out in parameters  $\mathbf{W}$ ,  $\mathbf{S}$ ,  $\theta$  among a batch of training samples  $(\mathbf{x}_{[i]}, \alpha_{[i]}^*)$ ,  $i = 1, \dots, M$  with stochastic gradient descent and backpropagation through time:

$$\mathbf{L}(\mathbf{W}, \mathbf{S}, \theta) = \frac{1}{2} \sum_{i=1}^M \|\alpha_{[i]}^* - \alpha_{[i]}^{(\mathbf{K})}\|_2^2. \quad (5)$$

### Advantages:

- Matrix multiplications and soft-thresholding are *fast*.
- The algorithm has an *adjustable iteration count  $\mathbf{K}$* .
- The method is *sparsity adaptive*, as  $\theta$  is learnable.
- The *problem is reduced*, as only a correlated subset of all possible signals and sparse codes appear in an actual dataset.

### Limitations:

- The scheme is limited to the case of the *linear and unstructured* LASSO problem.

## Hypotheses

### Our hypotheses:

- **LISTA may generalize to the structured case** via replacing outputs  $\alpha^*$  with binary group activations  $\left( \frac{\|\alpha_{\mathbf{G}}^*\|_2}{\sqrt{|\mathbf{G}|}} \right)_{\mathbf{G} \in \mathcal{G}} > \mathbf{0}$ , and computing a reduced pseudo-inverse.
- **LISTA may generalize to the kernelized case** via inputs  $\mathbf{k}(\mathbf{D}, \mathbf{x})$  or  $\left( \frac{\|\mathbf{k}(\mathbf{D}_{\mathbf{G}}, \mathbf{x})\|_2}{\sqrt{|\mathbf{G}|}} \right)_{\mathbf{G} \in \mathcal{G}}$ .
- **LISTA may bypass kernel computations** for further speedup via inputs  $\mathbf{x}$  (as long as they are vectorizable).

To test our hypotheses, we implemented KGLASSO in Matlab and LISTA in Theano.

## Experimental Setup

For our numerical experiment, we used *3D position features of uppercase air-handwriting characters* from the **6D Motion Gesture Database (6DMG)** [5]. The data contained **26** characters (A to Z) each repeated at most **10** times by **25** subjects. We normalized the data and uniformly interpolated each sample to length **128**. We then partitioned the set into train (**17** subjects) and test (**8** subjects) parts. Train samples were averaged after fixing both the subject and the character: the resulting **442** mean curves served as the dictionary  $\mathbf{D}$ . The test set supplied **102336**  $\mathbf{x}_i$  signals by generating random convex combinations with additive noise after fixing both the subject and the character. Group structure  $\mathcal{G}$  was induced by subjects ( $\mathbf{G}$ ).  $\mathbf{k}(\mathbf{x}, \mathbf{y})$  was chosen to be the **Global Alignment time series kernel** [6] with parameter  $\sigma = \mathbf{0.9}$ . We then set  $\lambda = \mathbf{0.001}$  and computed ground truth binary group activations  $\left( \frac{\|\alpha_{\mathbf{G}}^*\|_2}{\sqrt{|\mathbf{G}|}} \right)_{\mathbf{G} \in \mathcal{G}} > \mathbf{0}$  according to equation (2). The average number of active groups thus became **5.2654**. The LISTA matrix  $\mathbf{W}$  was pretrained with tied weights. For measuring **multi-label classification** performance, we applied an **80%-10%-10%** training-validation-testing shuffle-and-split scheme.

## Results

Micro-averaged **multi-label classification results** were as follows.

Test set performance metric	$\mathbf{k}(\mathbf{D}, \mathbf{x})$		$\left( \frac{\ \mathbf{k}(\mathbf{D}_{\mathbf{G}}, \mathbf{x})\ _2}{\sqrt{ \mathbf{G} }} \right)_{\mathbf{G} \in \mathcal{G}}$		$\mathbf{x}$	
	$\mathbf{K} = 1$	$\mathbf{K} = 2$	$\mathbf{K} = 1$	$\mathbf{K} = 2$	$\mathbf{K} = 1$	$\mathbf{K} = 2$
Loss function value	<b>0.300</b>	<b>0.329</b>	<b>1.497</b>	<b>1.077</b>	<b>1.323</b>	<b>1.052</b>
Accuracy	<b>0.962</b>	<b>0.961</b>	<b>0.793</b>	<b>0.850</b>	<b>0.758</b>	<b>0.820</b>
Precision	<b>0.930</b>	<b>0.924</b>	<b>0.824</b>	<b>0.781</b>	<b>0.599</b>	<b>0.682</b>
Recall	<b>0.950</b>	<b>0.954</b>	<b>0.418</b>	<b>0.712</b>	<b>0.671</b>	<b>0.789</b>
$\mathbf{F}_1$ score	<b>0.940</b>	<b>0.939</b>	<b>0.555</b>	<b>0.745</b>	<b>0.633</b>	<b>0.732</b>

## Conclusion

### Our findings:

- **LISTA can generalize to the structured and kernelized KGLASSO case.**
- **Single layer is already capable.**
- **Mapping from kernels  $\mathbf{k}(\mathbf{D}, \mathbf{x})$  is very accurate.**
- **Mapping from signals  $\mathbf{x}$  the performance is considerable.**

Results may improve for larger  $\mathbf{K}$  and with Convolutional neural network (CNN).

## References

- [1] László A Jeni, András Lőrincz, Zoltán Szabó, Jeffrey F Cohn, and Takeo Kanade. Spatio-temporal event classification using time-series kernel based structured sparsity. In *Computer Vision–ECCV 2014*, pages 135–150. Springer, 2014.
- [2] Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, 2(1):183–202, 2009.
- [3] Karol Gregor and Yann LeCun. Learning fast approximations of sparse coding. In *Proc. of the 27th Int. Conf. on Machine Learning (ICML-10)*, pages 399–406, 2010.
- [4] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996.
- [5] Mingyu Chen, Ghassan AlRegib, and Bing-Hwang Juang. 6dmg: A new 6d motion gesture database. In *Proc. of the 3rd Multimed. Sys. Conf.*, pages 83–88. ACM, 2012.
- [6] Marco Cuturi. Fast global alignment kernels. In *Proc. of the 28th Int. Conf. on Machine Learning (ICML-11)*, pages 929–936, 2011.