LEARNING DATA TRIAGE: LINEAR DECODING WORKS FOR COMPRESSIVE MRI

Yen-Huan Li and Volkan Cevher

Laboratory for Information and Inference Systems École Polytechnique Fédérale de Lausanne

ABSTRACT

The standard approach to compressive sampling considers recovering an unknown deterministic signal with certain known structure, and designing the sub-sampling pattern and recovery algorithm based on the known structure. This approach requires looking for a good representation that reveals the signal structure, and solving a non-smooth convex minimization problem (e.g., basis pursuit). In this paper, another approach is considered: We learn a good sub-sampling pattern based on available training signals, without knowing the signal structure in advance, and reconstruct an accordingly sub-sampled signal by computationally much cheaper linear reconstruction. We provide a theoretical guarantee on the recovery error, and show via experiments on real-world MRI data the effectiveness of the proposed compressive MRI scheme.

Index Terms— Compressive sampling, magnetic resonance imaging (MRI), least squares estimation, sub-modular minimization, statistical learning

1. INTRODUCTION

The standard theory of compressive sampling (CS) considers recovering an unknown deterministic signal with certain *known* structure, and designing sampling and recovery schemes based on the known structure [11]. For example, if the unknown signal is known to be sparse, one can measure it by a sub-sampling matrix satisfying the restricted isometry property (RIP), and apply basis pursuit to obtain an estimate of the signal [5, 6]. Similar ideas can be extended to low-rank matrix recovery [4], and in general signal recovery problems where the signal structures can be encoded by atomic norms or other convex functions [2, 8, 10].

Despite its success in many applications, we note that there are some undesired features of the standard CS theory:

1. The signal structure must be known *in advance*. This usually requires seeking for a good signal representation to reveal the signal structure, a non-trivial task called dictionary learning [17].

2. The recovery scheme is *computationally expensive*. Typical examples are basis pursuit and the Lasso, which are both non-smooth convex optimization problems.

While those features seem to be necessary according to existing literature on CS, in some applications, the real-world setting can deviate from the standard setting of CS. This creates an opportunity of getting rid of those undesired features.

We focus on one important observation which the standard CS theory does not take into consideration—we usually have training signals, i.e., signals that are given and similar to the unknown signal in some sense.

In fact, practitioners are indeed applying this learningbased approach in a naïve way. For example, it is by examining a large amount of real-world images that we discovered sparsity or more sophisticated structures, under proper representations [3, 7, 15]. Although this naïve learning procedure can be made rigorous and automated by dictionary learning, training signals are still required.

In this paper, we propose alternative to compressive sampling, and apply it to compressive magnetic resonance imaging (MRI). The proposed scheme *automatically* adapts to the given training signals, without any *a priori* knowledge on the signal structure. We highlight the following contributions:

- 1. We propose a novel statistical learning view point to the compressive sampling problem, which allows us to study the effect of training signals.
- 2. Our compressive MRI scheme is computationally efficient: The learning procedure can be cast as a combinatorial optimization problem, which can be exactly solved by an efficient algorithm; the recovery algorithm we consider is simply least-squares (LS) reconstruction.
- 3. In contrast to the standard approach using random subsampling patterns [6, 14, 16], our sub-sampling scheme is fixed given the training signals, and hence simpler for implementation.
- 4. We provide a theoretical guarantee on the reconstruction error, and characterize its dependence on the number of training signals.

This work was supported in part by ERC Future Proof, SNF 200021-146750 and SNF CRSII2-147633.

5. We show via experiments on real MRI images that the reconstruction error performance of the proposed scheme is comparable to the performance using a finely-tuned sub-sampling pattern given in [14].

2. REVIEW OF EXISTING APPROACHES

Compressive MRI is essentially a linear inverse problem. The goal is to recover an unknown signal $x^{\natural} \in \mathbb{C}^{p}$, given a subsampling pattern $\Omega \subset \{1, \ldots, p\}$ with $|\Omega| = n$ for some n < p, and the outcome of compressive sampling:

$$y := P_{\Omega} \mathcal{F} x^{\sharp}$$

where $\mathcal{F}: \mathbb{C}^p \to \mathbb{C}^p$ is the Fourier transform matrix, and $P_{\Omega}: \mathbb{C}^p \to \mathbb{C}^n$ is a linear operator that only keeps entries of $\mathcal{F}x^{\natural}$ indexed by Ω . In practice, x^{\natural} is usually a 2D or 3D object, and \mathcal{F} should be replaced by the corresponding multidimensional Fourier transform.

Existing approaches to compressive MRI can be briefly summarized as follows:

- Find a wavelet transform matrix Ψ : C^p → C^p, such that x^β = Ψ⁻¹z^β and z^β possesses certain *structure*. For example, the sparsity of z^β and smoothness of x^β were exploited in [14], the tree sparsity of z^β was considered in [9], and the multi-level sparsity of z^β was considered in [16].
- Choose a *random* sub-sampling pattern Ω and sample *Fx^{\\phi}* accordingly; the probability distribution might be dependent on the knowledge about the structure of z^{\\phi} [14, 16].
- Finally, apply *non-linear* decoding algorithms to reconstruct x^k. The standard basis pursuit estimator was considered in [16]. A basis pursuit like estimator minimizing a linear combination of the l₁-norm and the total variation semi-norm was proposed in [14]. A closely-related Lasso like estimator with the l₁-norm and total variation semi-norm penalization was considered in [19]. A similar Lasso like estimator with one additional penalization term for tree sparsity was introduced by [9].

We note that existing approaches essentially follow the standard theory of compressive sampling, and hence inherit the two undesired features which we mentioned in the introduction.

3. LEARNING DATA TRIAGE

The standard approach to compressive MRI models x^{\natural} as a deterministic unknown signal. Here we adopt another modeling philosophy: We assume that x^{\natural} is a random vector following some *unknown* probability distribution \mathbb{P} , and we have

access to *m* training signals $x_1, \ldots, x_m \in \mathbb{C}^p$, which are independent and identically distributed random vectors also following \mathbb{P} , and are independent of x^{\natural} . Note that this is different from Bayesian compressive sampling [13], as \mathbb{P} is unknown in our model.

We consider LS reconstruction. For any given subsampling pattern Ω , the estimator has an explicit form:

$$\hat{x}_{\Omega} = \arg\min_{x} \left\{ \|y - P_{\Omega} \mathcal{F}x\|_{2}^{2} : x \in \mathbb{R}^{p} \right\}$$
$$= \mathcal{F}^{H} P_{\Omega}^{T} y.$$

Once the reconstruction scheme is fixed, the only issue is to choose Ω that optimizes the resulting estimation performance.

We show in Section 8.1 that for any given Ω , the expected normalized reconstruction error satisfies

$$\mathbb{E} \frac{\left\|\hat{x}_{\Omega} - x^{\natural}\right\|_{2}^{2}}{\left\|x^{\natural}\right\|_{2}^{2}} = 1 - \mathbb{E} f_{\Omega}(x), \tag{1}$$

where the expectations are with respect to $x^{\natural} \sim \mathbb{P}$ and $x \sim \mathbb{P}$, respectively, and we define

$$f_{\Omega}(x) := \frac{\|P_{\Omega}\mathcal{F}x\|_{2}^{2}}{\|x\|_{2}^{2}}$$

for convenience. This implies that the optimal sub-sampling pattern Ω , denoted by $\Omega_{opt.}$, is given by any solution of the following combinatorial optimization problem:

$$\Omega_{\text{opt}} \in \arg\max_{\Omega} \left\{ \mathbb{E} f_{\Omega}(x) : \Omega \subset \{1, \dots, p\}, |\Omega| = n \right\}.$$
(2)

However, since \mathbb{P} is assumed unknown, the optimization problem is not tractable.

Motivated by the idea of empirical risk minimization in statistical learning theory [18], we make use of the training signals and approximate $\Omega_{opt.}$ via any solution of the optimization problem:

$$\Omega_m \in \arg\max_{\Omega} \left\{ \hat{\mathbb{E}}_m f_{\Omega}(x) : \Omega \subset \{1, \dots, p\}, |\Omega| = n \right\}$$
(3)

where $\hat{\mathbb{E}}_m$ denotes the expectation with respect to the empirical measure, i.e.,

$$\hat{\mathbb{E}}_m f_{\Omega}(x) := \frac{1}{m} \sum_{i=1}^m \frac{\|P_{\Omega} \mathcal{F}_x_i\|_2^2}{\|x_i\|_2^2}$$

This optimization problem is tractable, because we only need to solve it for any realization of the training signals x_1, \ldots, x_m . Note that then $\hat{\mathbb{E}}_m f_{\Omega}(x)$ depends on x_1, \ldots, x_m and is random, and so does Ω_m .

The overall systems is summarized as follows:

1. Find a sub-sampling pattern Ω_m by (3).

2. Sub-sample x^{\natural} using Ω_m and obtain the measurement outcome

$$y := P_{\Omega_m} \mathcal{F} x^{\natural}$$

3. Recover x^{\natural} by

$$\hat{x} := \mathcal{F}^H P_{\Omega_m}^T y.$$

On Computing Ω_m : The optimization problem (3) is modular, and hence can be exactly solved by a simple greedy algorithm [12]: Let ϕ_i^T be the *i*-th row of Φ . Compute the values

$$v_i := \frac{1}{m} \sum_{i=1}^m \left(\phi_i^T x_i \right)^2,$$

and set Ω_m as the set of indices corresponding to the largest $n v_i$'s. The computational complexity is dominated by computation of v_i 's, which behaves as $\mathcal{O}(mp^2)$ in general, and $\mathcal{O}(mp \log p)$ if Φ is suitably structured, such as the Fourier and Hadamard matrices.

4. PERFORMANCE ANALYSIS

We analyze the reconstruction error of the proposed learningbased compressive sampling system.

If we could solve the optimization problem (2), the estimation performance would be given by

$$\mathbb{E} \frac{\left\|\hat{x}_{\Omega_{\text{opt}}} - x^{\natural}\right\|_{2}^{2}}{\left\|x^{\natural}\right\|_{2}^{2}} = 1 - \varepsilon_{\mathbb{P}}$$

where we define

$$\varepsilon_{\mathbb{P}} := \max_{\Omega} \left\{ \mathbb{E} f_{\Omega}(x) : \Omega \subset \{1, \dots, p\}, |\Omega| = n \right\}.$$

Note that $\varepsilon_{\mathbb{P}}$ is a constant for any given \mathbb{P} , independent of the training signals.

Since now the optimization problem (2) is replaced by its empirical version (3), a reasonable guess is that the estimation performance would behave as

$$\mathbb{E} \frac{\left\|\hat{x}_{\Omega_m} - x^{\natural}\right\|_2^2}{\left\|x^{\natural}\right\|_2^2} \le 1 - \varepsilon_{\mathbb{P}} + \varepsilon_m,$$

for some $\varepsilon_m > 0$ with high probability (with respect to the training signals x_1, \ldots, x_m), and ε_m should converge to 0 as $m \to \infty$. The following proposition verifies this guess.

Proposition 4.1. *For any* $\beta \in (0, 1)$ *, we have*

$$\varepsilon_m \le \sqrt{\frac{2}{m}} \left[\log \binom{p}{n} + \log \frac{2}{\beta} \right]$$

with probability at least $1 - \beta$.

This means a size of training signals of the order $\mathcal{O}(n \log p)$ suffices to have small enough ε_m , with high probability. We note that this is a worst-case guarantee, as it is *distribution-independent*. In practice, *m* can be much smaller.



Fig. 1. First row: the subsampling maps of the tuned random variable sampling scheme [14]. Second row: the maps given by our learning-based approach.

5. NUMERICAL RESULTS

We use a 3-dimensional dataset of raw knee-images data given in k-space.¹ We first take an inverse Fourier transform along the z-axis and eliminate low energy the z-slices that are close to the boundary of the datacube. These are noise-like slices that do not exhibit any knee feature as they are close to the skin of the patient. We then investigate subsampling schemes in the $320 \times 320 x - y$ Fourier plane, which corresponds to compressive sampling for each z-slice.

We pick the first 10 of the patients in the given dataset for training and test the learned subsampling maps on the remaining 10 patients. We compare our learning based approach to the variable density function proposed by [14], which is parametrized by the radius of fully sampled region, r, and the polynomial degree, d. We tune the values of r and d so that they yield the highest average PSNR on the training data.

Figure 1 illustrates the best performing randomized indices and our learned set of indices in the x - y plane of the k-space. Both the variable density approach [14] and our learning-based approach concentrates its sampling budget on the low frequencies, however the latter is endowed with the capability to adapt its frequency selection to the frequency content of the training signals instead of assuming a circularly symmetric selection.

 Table 1. Average PSNR on the test data

Indices	Sampling rate		
	6.25%	12.50%	25%
Best-n approx.	$25.29~\mathrm{dB}$	26.36 dB	28.35 dB
Lustig et al.	24.51 dB	25.11 dB	26.05 dB
This work	24.66 dB	25.18 dB	26.12 dB

¹Available at http://mridata.org/fullysampled



Fig. 2. MRI reconstructions of both schemes at different subsampling rates for a knee slice of patient #13, whose fully sampled reconstruction is shown on the top left.

Table 1 shows the performance of both approaches on the test data, in addition to the error lower-bounds obtained by the best n-sample approximations with respect to the Fourier basis. It appears that the learning based approach slightly outperforms the randomized variable density based approach.

However, the slight numerical improvements are actually accentuated when we look at the details of reconstructions, shown in Figure 2 for the test Patient #13. It is clear that the learning-based reconstructions provide more details especially for 6.25% and 12.5%.

6. DISCUSSIONS

The essential idea of the learning-based approach can be summarized as follows: Fix a decoder, and find the optimal subsampling pattern that minimizes the corresponding expected recovery error, which can be approximated by empirical risk minimization. The performance is essentially determined by the distribution of signal ensemble.

In this paper, we consider the linear decoder for computational efficiency, and it works well on the ensemble of MRI images. For other signal ensembles, it is possible to have a better recovery error performance by a non-linear decoder, such as basis pursuit or the Lasso, and realize a trade-off between computational complexity and recovery performance. Note that the idea of the learning-based approach still applies, while the empirical risk minimization formulation for choosing the sub-sampling pattern should be modified accordingly given the decoder. We are currently working in this research direction.

7. ACKNOWLEDGEMENT

The authors would like to thank Baran Gözcü and Luca Baldassarre for providing numerical results, and Jonathan Scarlett for the complexity discussion.

8. PROOFS

8.1. Proof of (1)

. 2

In fact, the equality holds deterministically, as

$$\begin{split} \left\| \hat{x}_{\Omega} - x^{\sharp} \right\|_{2}^{2} \\ &= \left\| \hat{x}_{\Omega} \right\|_{2}^{2} - 2 \left\langle \hat{x}_{\Omega}, x^{\natural} \right\rangle + \left\| x^{\natural} \right\|_{2}^{2} \\ &= \left\| \mathcal{F}^{H} P_{\Omega}^{T} P_{\Omega} \mathcal{F} x^{\natural} \right\|_{2}^{2} - 2 \left\langle F^{H} P_{\Omega}^{T} P_{\Omega} \mathcal{F} x^{\natural}, x^{\natural} \right\rangle + \left\| x^{\natural} \right\|_{2}^{2} \\ &= \left\| P_{\Omega} \mathcal{F} x^{\natural} \right\|_{2}^{2} - 2 \left\| P_{\Omega} \mathcal{F} x^{\natural} \right\|_{2}^{2} + \left\| x^{\natural} \right\|_{2}^{2}. \end{split}$$

In the third equality, we used the fact that $AA^{\dagger}A = A$ for any matrix A and its Moore-Penrose generalized inverse A^{\dagger} , by setting $A := P_{\Omega} \mathcal{F}$.

8.2. Proof of Proposition 4.1

It suffices to choose ε_m such that with probability at least $1 - \beta$,

$$\Delta_m := \mathbb{E} f_{\Omega_{\text{opt}}}(x) - \mathbb{E} f_{\Omega_m}(x) \le \varepsilon_m.$$

We note that

$$\begin{split} \Delta_m &= \left(\mathbb{E} f_{\Omega_{\text{opt}}}(x) - \hat{\mathbb{E}}_m f_{\Omega_{\text{opt}}}(x) \right) + \\ \left(\hat{\mathbb{E}}_m f_{\Omega_{\text{opt}}}(x) - \hat{\mathbb{E}}_m f_{\Omega_m}(x) \right) + \\ \left(\hat{\mathbb{E}}_m f_{\Omega_m}(x) - \mathbb{E} f_{\Omega_m}(x) \right). \end{split}$$

The second summand on the right-hand side cannot be positive by definition. Then we have

$$\Delta_m \le 2 \max_{\Omega} \left\{ \left| \hat{\mathbb{E}}_m f_{\Omega}(x) - \mathbb{E} f_{\Omega}(x) \right| : \Omega \in \mathcal{A} \right\}$$

where $\mathcal{A} := \{\Omega : \Omega \subset \{1, \dots, p\}, |\Omega| = n\}.$

As the random variables $f_{\Omega}(x)$ are bounded $(0 \le f_{\Omega}(x) \le 1)$, we can use Hoeffding's inequality and the union bound to obtain an upper bound of Δ_m that holds with high probability, as in [1, B.3].

9. REFERENCES

- J.-Y. Audibert and O. Bousquet, "Combining PAC-Bayesian and generic chaining bounds," *J. Mach. Learn. Res.*, vol. 8, pp. 863–889, 2007.
- [2] F. Bach, "Learning with submodular functions: A convex optimization perspective," *Found. Trends Mach. Learn.*, vol. 6, no. 2–3, pp. 145–373, 2013.

- [3] R. G. Baraniuk, V. Cevher, M. F. Duarte, and C. Hegde, "Model-based compressive sensing," *IEEE Trans. Inf. Theory*, vol. 56, no. 4, pp. 1982–2001, Apr. 2010.
- [4] E. Candès and B. Recht, "Exact matrix completion via convex optimization," *Found. Comput. Math.*, vol. 9, pp. 717–772, 2009.
- [5] E. J. Candès, "The restricted isometry property and its implications for compressed sensing," C. R. Acad. Sci. Paris, Ser. I, vol. 346, pp. 589–592, 2008.
- [6] E. J. Candès, J. Romberg, and T. Tao, "Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information," *IEEE Trans. Inf. Theory*, vol. 52, no. 2, pp. 489–509, Feb. 2006.
- [7] V. Cevher, "Learning with compressible priors," in Adv. Neural Information Processing Systems, 2009.
- [8] V. Chandrasekaran, B. Recht, P. A. Parrilo, and A. S. Willsky, "The convex geometry of linear inverse problems," *Found. Comput. Math.*, vol. 12, pp. 805–849, 2012.
- [9] C. Chen and J. Huang, "Compressive sensing MRI with wavelet tree sparsity," in Adv. Neural Information Processing Systems 25, 2012.
- [10] M. El Halabi and V. Cevher, "A totally unimodular view of structured sparsity," in 18th Int. Conf. Artificial Intelligence and Statistics, 2015.
- [11] S. Foucart and H. Rauhut, A Mathematical Introduction to Compressive Sensing. Basel: Birkhäuser, 2013.
- [12] S. Fujishige, Submodular Functions and Optimization, 2nd ed. Amsterdam: Elsevier, 2005.
- [13] S. Ji, Y. Xue, and L. Carin, "Bayesian compressive sensing," *IEEE Trans. Sig. Process.*, vol. 56, no. 6, pp. 2346– 2356, 2008.
- [14] M. Lustig, D. Donoho, and J. M. Pauly, "Sparse MRI: The application of compressed sensing for rapid MR imaging," *Magn. Reson. Med.*, vol. 58, pp. 1182–1195, 2007.
- [15] S. Mallat, A Wavelet Tour of Signal Processing: The Sparse Way, 3rd ed. Burlington, MA: Academic Press, 2009.
- [16] B. Roman, B. Adcock, and A. Hansen, "On asymptotic structure in compressed sensing," 2014, arXiv:1406.4178v2 [math.FA].
- [17] I. Tošić and P. Frossard, "Dictionary learning," *IEEE Sig. Process. Mag.*, vol. 28, no. 2, pp. 27–38, 2011.

- [18] V. N. Vapnik, "An overview of statistical learning theory," *IEEE Trans. Inf. Theory*, vol. 10, no. 5, pp. 988– 999, Sep. 1999.
- [19] J. Yang, Y. Zhang, and W. Yin, "A fast alternating direction method for TVL1-L2 signal reconstruction from partial Fourier data," *IEEE J. Sel. Topics Signal Process.*, vol. 4, no. 2, pp. 288–297, 2010.