Associative Memory using Dictionary Learning and Expander Decoding

Arya Mazumdar[†] and Ankit Singh Rawat*[‡]

[†]College of Information & Computer Science, University of Massachusetts Amherst, MA, USA [‡]Research Laboratory of Electronics, Massachusetts Institute of Technology, MA, USA E-mail: arya@cs.umass.edu, asrawat@mit.edu

November 30, 2016

Abstract

An associative memory is a framework of content-addressable memory that stores a collection of message vectors (or a *dataset*) over a neural network while enabling a neurally feasible mechanism to recover any message in the dataset from its noisy version. Designing an associative memory requires addressing two main tasks: 1) *learning phase*: given a dataset, learn a concise representation of the dataset in the form of a graphical model (or a neural network), 2) *recall phase*: given a noisy version of a message vector from the dataset, output the correct message vector via a neurally feasible algorithm over the network learnt during the learning phase. This paper studies the problem of designing a class of neural associative memories which learns a network representation for a large dataset that ensures correction against a large number of adversarial errors during the recall phase. Specifically, the associative memories designed in this paper can store dataset containing $\exp(n)$ n-length message vectors over a network with O(n) nodes and can tolerate $\Omega(\frac{n}{\operatorname{polylog} n})$ adversarial errors. This paper carries out this memory design by mapping the learning phase and recall phase to the tasks of dictionary learning with a square dictionary and iterative error correction in an expander code, respectively.

1 Introduction

Associative memories aim to address a problem that naturally arises in many information processing systems: given a dataset \mathcal{M} which consists of n-length vectors, design a mechanism to concisely store this dataset so that any future query corresponding to a noisy version of one of the vectors in the dataset can be mapped to the correct vector. An associative memory based solution to this problem is broadly required to have two key components: 1) dataset must be stored in the form of a neural network (graph) and 2) the mechanism to map a noisy query to the associated valid vector should be implementable in an iterative neurally feasible manner over the network (a *neurally feasible* algorithm employs only local computations at the nodes of the corresponding network based on the information obtained from their neighboring nodes). The tasks of learning the graph representation from the dataset and mapping erroneous vectors to the associated correct vectors are referred to as *learning phase* and *recall phase*, respectively.

The overarching goal of designing an associative memory that can store a large dataset (ideally containing $\exp(n)$ message vectors using a neural network with O(n) nodes) while ensuring robustness to a large number of errors (ideally $\Omega(n)$ errors) during the recall phase has led to multiple research efforts in the

^{*}This work was done when the author was with the Computer Science Department, Carnegie Mellon University, PA, USA.

literature. The binary Hopfield networks, as studied in [13,22], provide one of the earliest designs for the associative memories. Given a dataset containing binary vectors from $\{\pm 1\}^n$, Hopfield networks learn this dataset in the form of an n-node weighted graph by employing Hebbian learning [11], i.e., the weighted adjacency matrix of the graph is defined by summing the outer products of all message vectors in the dataset. However, in their most general form, these networks suffer from small capacity. In [22], McEliece et al. show that these networks can only store $O\left(\frac{n}{\log n}\right)$ message vectors when these messages correspond to arbitrary n-length binary vectors and the recall phase is required to tolerate linear $\Omega(n)$ random errors. This has motivated the researchers to look at various generalizations of Hopfield networks (see, [9,16,21,23,31] and references therein). However, these solutions again fail to simultaneously achieve both large capacity and error tolerance.

One remedy to small capacity is to design associative memories with structural assumptions on the dataset. This approach has been considered in [8, 12, 18–20, 27]. In particular in [8], Gripon et al. store a dataset comprising $O(n^2)$ sparse vectors in the form of cliques in a neural network. In [12], Hillar and Tran design a Hopfield network with n nodes that can store $\sim 2^{\sqrt{2n}}/n^{1/4}$ message vectors and is robust against n/2 random errors. In [18–20, 27], the message vectors that need to be stored are assumed to constitute a subspace. In [18, 19, 27], the task is to learn a bipartite factor graph of the linear constraints satisfied by the dataset subspace. The error correction during recall phase is then performed by running a belief propagation algorithm [26] over the bipartite graph. In [18], Karbasi et al. work with a model where the message vectors in the dataset have overlapping sets of coordinates so that shortened vectors obtained by restricting the original message vectors to each of these overlapping sets belong to a subspace. Under this model, they design associative memories that can store exponential number (in n) of message vectors while correcting linear number (in n) of random errors during the recall phase.

The results in [18] hinge on the fact that the learning phase of their memory design recovers a bipartite graph which has certain desirable structural properties that are required for belief propagation type decoders to converge. However, no guarantee of recovering such a bipartite graph during the learning phase is provided in [18] even when we assume the subspace associated with the dataset has one such graphical representation to begin with. Recognizing the requirement of learning correct bipartite graph during the learning phase, Mazumdar and Rawat explore a sparse recovery based approach to design associative memories with the subspace dataset model in [20]. This approach assumes that the dataset belongs to a subspace whose orthogonal subspace has null space property, a sufficient condition for sparse signal recovery. This allows one to learn any basis for the orthogonal subspace during the learning phase and then recast the recall phase as a sparse recovery problem [4,6]. The approach in [20] also allows for the strong error model containing adversarial errors. Specifically, [20] considers two candidate signal models which contain n-length message vectors and utilize O(n) sized neural networks to store the signals. The two models have the datasets of sizes $\exp(n^{3/4})$ and $\exp(r)$ with $1 \le r \le n$, respectively. Furthermore, the designed associative memories based on these two signal models respectively allow for recovery from $\Omega(n^{1/4})$ and $\Omega\left(\frac{n-r}{\log^6 n}\right)$ adversarial errors in a neurally feasible manner.

In this paper, we also follow the subspace model as in [18–20, 27]. We assume the dataset to form a subspace which is defined by *sparse linear constraints*. The model of sparse linear constraints are quite natural and less restrictive than the previous models of works such as [20]. Note that this signal model is similar to the model explored in Karbasi et al. [18]. However, our approach and contributions differ from [18], as we ensure that the learning phase *provably* generates the correct bipartite graph which can guarantee the error correction from a large number of errors using an iterative algorithm during the recall phase. We also note that similar to [20] we work with the stronger error model involving adversarial errors, but our scheme is superior to that of [20] in terms of storage capacity (see, Theorems 1, 3) and number of correctable adversarial errors (improvement by poly-log factors, see, Theorems 1, 2, 3). We want to point out that the main technical challenge in associative memory is not to individually design the learning or

recall phases, but to interface them in a way that is consistent with the operations of both phases, and to give an end-to-end performance guarantee.

Here, we note that the problem of designing an associative memory is closely related to the well studied nearest neighbor search (NNS) problem and its relaxation approximate nearest neighbor search (A-NNS) problem [2, 14, 28, 32]. The solutions to the A-NNS problem enable one to store a dataset in such a manner that noisy versions of the vectors in a dataset (with bounded noise) can be mapped to the correct vectors. Additionally, the A-NNS solutions do not put assumptions on the dataset. However, this comes at the cost of removing the requirement of having a fast iterative or neurally feasible recall phase. Furthermore, the A-NNS solutions, especially based on locally sensitive hashing [10, 14] have large space complexity, i.e., polynomial in size of dataset. We note that the A-NNS solutions are very much aligned to the vector (image) retrieval task [7, 17, 33] which need not have a neurally feasible retrieval algorithm.

The rest of the paper is organized as follows. In Sec. 2, we define the dataset model considered in this paper and present the main results of this paper along with key techniques and ideas involved in establishing those results. Sec. 3 is dedicated to the proof of the main theorem. In Sec. 3.1, we describe the learning phase of the associative memory design results along with the relevant technical details. In Sec. 3.2, we present an iterative error correction algorithm which is employed during the recall phase of the designed associative memory. This analysis of the algorithm relies on the expansion properties of the bipartite graph which defines the dataset and is learnt during the learning phase. We conclude the paper with some comments on performance in Sec. 5.

2 Main results and techniques

2.1 Model for datasets

We focus on the associative memories based on the operations on \mathbb{R} , the set of real numbers. In our first model, we consider the message patterns to be vectors over \mathbb{R} . In the second model we comment on neural associative memories storing binary message patterns that are obtained by our approach.

2.1.1 Dataset over real numbers: the sparse-sub-Gaussian model

We assume the message set to form a linear subspace defined by sparse linear constraints over \mathbb{R} . Let $\mathcal{M} \subseteq \mathbb{R}^n$ denote the set of message vectors (signals) that need to be stored on the associative memory. Let B be an $m \times n$ matrix comprising the linear constraints that define the message set \mathcal{M} . In particular, we have

$$Bx = 0 \quad \forall x = (x_1, x_2, \dots, x_n) \in \mathcal{M}. \tag{1}$$

In order to fully specify the message set \mathcal{M} , we still need to provide a stochastic model for the matrix B. Towards this, we consider a random ensemble of sparse matrices. For each $j \in [n] := \{1, 2, \dots, n\}$, we consider the following experiment. We pick d element uniformly at random with replacement from the set [m]. Let \mathcal{N}_j denote the set comprising these randomly picked elements. For $1 \le i \le m$, $1 \le j \le n$, we define

$$\xi_{i,j} = \begin{cases} 1 & \text{if } i \in \mathcal{N}_j \subset [m] \\ 0 & \text{otherwise.} \end{cases}$$
 (2)

Let $\{R_{i,j}\}_{1 \leq i \leq m, \ 1 \leq j \leq n}$ be a collection of independent and identically distributed (i.i.d.) sub-Gaussian random variables. Given the random variables, $\{\xi_{i,j}, R_{i,j}\}_{1 \leq i \leq m, \ 1 \leq j \leq n}$, we assume that the (i,j)-th entry

of the matrix B is defined as

$$B_{i,j} = \xi_{i,j} R_{i,j} \in \mathbb{R} \quad \text{for } 1 \le i \le m, \ 1 \le j \le n. \tag{3}$$

Through out this paper, we refer to this model for the dataset to be stored on a neural associative memory as *sparse-sub-Gaussian model*. We work with various values of d which we specify while stating different parameters that we obtain for the designed associative memories in Sec. 2.2.

This model is a quite natural random model of bipartite graphs that allow for multi-edges. Indeed, consider a bipartite graph with disjoint sets of vertices [n] (variable nodes) and [m] (check nodes). There are d edges out of each variable node, being incident on uniformly and independently chosen vertices from the check nodes.

Remark 1. The requirement on $R_{i,j}$ is quite generic as it allows for many distributions. For example, we can assume that $R_{i,j}$ belongs to a finite set of integers $\{-L, -L+1, \ldots, -1, 1, \ldots, L-1, L\}$. Similarly, in another setup, $R_{i,j}$ can be assumed to be a Gaussian random variable.

2.1.2 Binary dataset

Our model of binary dataset is same as above except for the fact that 1) $\mathcal{M} \subseteq \{+1, -1\}^n$, and 2) $R_{i,j}$ is uniform over $\{+1, -1\}$ in (3). The condition of (1) must be satisfied for any $x \in \mathcal{M}$.

2.2 Our main results

We establish that, for a dataset \mathcal{M} corresponding to the null-space defined by the matrix B, the said matrix B can be exactly recovered from the dataset in polynomial time. Recall that there can be many sets of basis-vectors for the null-space of \mathcal{M} . Still, we claim that it is possible to accurately recover the matrix B that has been generated by the sparse-sub-Gaussian model described above.

It is essential for us that we recover the matrix B exactly. Being generated by the random model defined above, B exhibits certain graph expansion property that is necessary for our recall phase to be successful. This matrix B enables the error correction during the recall phase with the help of a simple iterative (neurally feasible algorithm). We summarize the parameters achieved by such memory as follows.

Theorem 1. Suppose that c, c', c'' > 0 are three constants. Let n be a large enough integer and $m = c \frac{n}{\log n}$. Assume that B is an $m \times n$ matrix generated from the sparse-sub-Gaussian model described in Sec. 2.1.1 with $c' \le d \le c'' \log n$, and $\mathcal{M} = \{x \in \mathbb{R}^n : Bx = 0\}$. Then, with high probability (w.h.p.) \mathcal{M} is an $n-m = n(1-c/\log n)$ dimensional subspace that can be stored in a neural network (learned in poly-time in the learning phase) while allowing for correct recovery from $\Omega(\frac{n}{d^2 \log^2 n})$ adversarial errors during the recall phase with a neurally feasible algorithm.

The proof of this theorem has been provided in Sec. 3. This result is obtained by utilizing a novel connection between recovering the matrix B defining the underlying dataset \mathcal{M} and the dictionary learning problem with a square dictionary as studied in [1,3,30]. Given access to the dataset \mathcal{M} , we can easily find a basis for the null-space of \mathcal{M} containing $m = n - \dim(\mathcal{M})$ n-length vectors. Let A denote the $m \times n$ matrix which has the m vectors in this basis as its rows. Note that the row vectors of B also span the subspace orthogonal to the dataset \mathcal{M} . Moreover, w.h.p., B is a full rank matrix. This implies that the following relationship holds w.h.p.,

$$A = DB, (4)$$

where D is an invertible $m \times m$ matrix. Note that recovering the matrix B from A is now equivalent to dictionary learning problem [24] where n columns of A and B corresponds to n observations and the associated coefficients, respectively. Furthermore the matrix D corresponds to a square dictionary [30].

As for the recall phase, we rely on the observations (as shown in Sec. 3.2) that w.h.p. the bipartite graph associated with the sparse random matrix B is an expander graph. Assume that we are given a noisy version y of a valid message vector $x \in \mathcal{M}$ such that we have

$$y = x + e \tag{5}$$

where e denotes the error vector. Recovering x from the observation y can be cast as a sparse recovery problem of recovering e from

$$z = By = B(x + e) = Be$$
.

If the bipartite graphs associated with B is an expander graph (which holds w.h.p.), we can solve this sparse recovery problem by an efficient and iterative algorithm [15] which is motivated by the decoding algorithm of expander codes [29] in coding theory literature.

Due to the sample complexity requirements for efficient square-dictionary learning algorithms [1,3,30], the above model allows us to store datasets that satisfy at most $O(\frac{n}{\log n})$ linear constraints. However if we allow for a learning-phase that takes quasi-polynomial time, then it is possible to store restricted datasets that satisfy $m = \Theta(n)$ sparse-linear constraints. We summarize the result below.

Theorem 2. Let n be a large enough integer and m=cn for a some constant c<1/200. For a large enough constant C>0, let B be an $m\times n$ matrix generated from the sparse-sub-Gaussian model described in Sec. 2.1.1 with $d=C\log n$ and $\mathcal{M}=\{x\in\mathbb{R}^n:Bx=0\}$. Then w.h.p., \mathcal{M} is an n-m=n(1-c) dimensional subspace that can be stored in a neural network (learned in quasi-polynomial-time in the learning phase) while allowing for error correction from $\Omega(\frac{n}{\log^2 n})$ adversarial errors during the recall phase with a neurally feasible algorithm.

While in terms of storage capacity this theorem is inferior to that of Theorem 1, it may represent some datasets better, and has better error correction capability. While the recall phase of this algorithm works same as above, for the learning phase we can no longer rely on the dictionary-learning algorithms. Instead we do an exhaustive search over all possible sparse vectors to find out a sparse basis for the null-space of $\mathcal M$ which end up taking a quasi-polynomial time, if we choose parameters suitable for the recall phase. We here crucially use the fact that for m=cn and $d=C\log n$ such a sparse basis is unique, which can be obtained from the results of [30]. The proof of the recall phase for this theorem remains same as that of Theorem 1.

Finally, while both Theorems 1 and 2 have their counterparts when storing binary vectors, we present only one result for brevity. A sketch of the proof of the following theorem has been given in Sec. 4.

Theorem 3 (Binary dataset). Suppose that c, c', c'' > 0 are three constants. Let n be a large enough integer such that $m = c \frac{n}{\log n}$. Assume that B is an $m \times n$ matrix generated from the binary dataset model described in Sec. 2.1.2 with $c' \le d \le c'' \log n$ and $\mathcal{M} = \{x \in \{\pm 1\}^n : Bx = 0\}$. Then w.h.p., $|\mathcal{M}| = \exp(n - \alpha n \log(d \log n)/\log n)$ for a constant α and \mathcal{M} can be stored in a neural network (learned in polynomial-time in the learning phase) while allowing for error correction from $O(\frac{n}{d^2 \log^2 n})$ adversarial errors during the recall phase with a neurally feasible algorithm.

3 Proof of Theorem 1

3.1 Learning phase of associative memory design

As discussed in the previous section, under the dataset model considered in this paper, the learning phase of the associative memory design can be mapped to the problem of dictionary learning with a square dictionary. The very same dictionary learning problem with slightly different random model for the coefficient vector has been studied in [1,3,30]. In Appendix A, we briefly describe this line of work along with the results that

are used in this paper. We then utilize the dictionary learning algorithm used in [1] to exactly learn the matrix B which define our dataset and comment on the modifications required in the analysis of Adamczak [1] to obtain guarantees on the performance of this algorithm.

3.1.1 Exact recovery of the matrix B

Our learning phase constitutes learning the matrix B exactly from the dataset \mathcal{M} . Utilizing the dictionary learning algorithm from [1], we design the learning phase for an associative memory storing the message set described in Sec. 2.1. The learning phase consists of the following two steps.

- 1. Given the message vectors from the dataset \mathcal{M} , first construct a basis for the subspace orthogonal to the dataset subspace $\mathcal{M} = \{x : Bx = 0\} \subset \mathbb{R}^n$ with $\dim(\mathcal{M}) = n m$.
- 2. Let $A \in \mathbb{R}^{m \times n}$ denote the basis obtained in the previous step. Since w.h.p. B is a full-rank matrix, we have

$$A = DB$$
,

where $D \in \mathbb{R}^{m \times m}$ is a non-singular matrix. Now employ the modified ERSpUD dictionary learning algorithm [1] with the matrix A as its input. Note that the algorithm outputs candidates for the matrices D and B. The method of this square-dictionary learning and the algorithm are summarized in Appendix A.

Next, we show that the proposed learning phase w.h.p. exactly recovers the matrix B. Note that the sparse-sub-Gaussian model used to generate B (cf. Sec. 2.1) slightly differs from the Bernoulli-sub-Gaussian model studied in [1, 30] (cf. Appendix A). In particular, for every $j \in [n]$, the distribution of the random variables $\{\xi_{i,j}: i \in [m]\}$ and $\{\eta_{i,j}: i \in [m]\}$ is different¹. However, this difference is not very crucial for the success of the learning algorithm as we still have independence among the random variables $\xi_{i,j}$ s which are indexed by different values of $j \in [n]$. We formalize the exact recovery guarantees for the matrix B in the following result.

Theorem 4. Let $B \in \mathbb{R}^{m \times n}$ be a matrix generated by the sparse-sub-Gaussian model (cf. Sec. 2.1) and \mathcal{M} be the associated dataset, i.e., $\mathcal{M} = \{x : Bx = 0\}$. Then there exists a constant c > 0 such that whenever we have $n \ge cm \log m$ the two step learning phase of the associative memory as described above exactly recovers the linear constraints in B with probability at least 1 - 1/n.

We refer the reader to Appendix A.1 for the proof of Theorem 4.

3.2 Recall phase of associative memory design

In this section we present an iterative algorithm which recovers the correct message vector among the dataset \mathcal{M} from its noisy version. The noisy observation is assumed to be corrupted at adversarially chosen coordinates. The correctness of the iterative algorithm relies on the observation that the bipartite graph associated with the matrix B which defines our dataset \mathcal{M} is a good expander graph. We first formalize this expansion property in the following result. We then present the iterative algorithm and show that it can provably tolerate $\Omega\left(\frac{n}{\text{polylog}n}\right)$ adversarial errors.

¹We focus on the sparse-sub-Gaussian model as opposed to the Bernoulli-sub-Gaussian model as the bipartite graph associated with the matrix B generated by the sparse-sub-Gaussian model is a good expander w.h.p. We utilize this fact while designing the recall phase for the proposed associative memory in Sec. 3.2.

3.2.1 Expansion property of the bipartite graph defined by B

Let $\mathcal{G}_B = (\mathcal{L} = [n], \mathcal{R} = [m], \mathcal{E}_B)$ be a bipartite graph where \mathcal{L} and \mathcal{R} denote the index sets of left and right vertices, respectively. The matrix B which defines our dataset \mathcal{M} gives the $m \times n$ adjacency matrix of the graph \mathcal{G} , i.e., for $\ell \in \mathcal{L}$ and $r \in \mathcal{R}$, we have an edge $(\ell, r) \in \mathcal{E}_B$ iff $B_{r,\ell} \neq 0$. More specifically, the weight of the edge $(\ell, r) \in \mathcal{E}_B$ is $w_{\ell,r} = B_{r,\ell}$. It follows from the sparse-sub-Gaussian model (cf. Sec. 2.1) which generates the random matrix B that every vertex in \mathcal{L} has degree d and each of the d neighbors for a vertex in \mathcal{L} are chosen uniformly at random from the set of right vertices \mathcal{R} with replacement. The following result states that expansion properties that hold for such a graph with high probability.

Proposition 1. Assume that $\epsilon > 0$ and $d = O(\frac{n}{m \log n})$. Let $\mathcal{G} = (\mathcal{L}, \mathcal{R}, \mathcal{E})$ be a random d-left regular graph where each of the d neighbors for a left vertex are chosen uniformly at random from the set of right vertices with replacement. Then, for a large enough n, w.h.p., \mathcal{G} is an $\left(\frac{m^2}{d^2n}, (1-\epsilon)d\right)$ -expander graph, where a bipartite graph is (t,l)-expander, if for every $\mathcal{S} \subseteq \mathcal{L}$ such that $|\mathcal{S}| \leq t$, we have $|\mathcal{N}(\mathcal{S})| \geq l|\mathcal{S}|$. Here, $\mathcal{N}(\mathcal{S}) \subseteq \mathcal{R}$ denotes the vertices in \mathcal{R} that are neighbors of vertices in \mathcal{S} .

Proof. Let's consider a set $\mathcal{S} \subseteq \mathcal{L}$ such that $|\mathcal{S}| = s \leq \frac{m^2}{d^2n}$. Let $\mathcal{T} \subseteq \mathcal{R}$ be a set of right vertices such that $|\mathcal{T}| < (1-\epsilon)ds$. The probability that $\mathcal{N}(\mathcal{S}) \subseteq \mathcal{T}$ is upper bounded by $\left(\frac{(1-\epsilon)ds}{m}\right)^{ds}$. Now, taking the union bound over all the sets $\mathcal{S} \subseteq \mathcal{L}$ such that $|\mathcal{S}| = s$ and the sets $\mathcal{T} \subseteq \mathcal{R}$ such that $|\mathcal{T}| < (1-\epsilon)ds$, the probability P_s that the the graph \mathcal{G} has a non-expanding set of size s, is upper bounded as follows.

$$P_{s} \leq {n \choose s} {m \choose (1-\epsilon)ds} ((1-\epsilon)ds/m)^{ds}$$

$$\leq e^{s+(1-\epsilon)ds} (n/s)^{s} ((1-\epsilon)ds/m)^{\epsilon ds}.$$
 (6)

We can rewrite (6) as,

$$P_s \le e^{s + (1 - \epsilon)ds} \left(\frac{dn}{m} \right)^s \left(\frac{ds}{m} \right)^{\epsilon ds - s}. \tag{7}$$

Now, using our assumption that $s \leq \frac{m^2}{d^2n}$, we obtain that

$$P_s \le e^{s + (1 - \epsilon)ds} \left(m/dn \right)^{\epsilon ds - 2s}. \tag{8}$$

Using union bound, we have that \mathcal{G} is not an $\left(\frac{m^2}{d^2n},(1-\epsilon)d\right)$ -expander with probability at most

$$\sum_{s=1}^{\frac{m^2}{d^2n}} P_s \le \frac{m^2}{d^2n} e^{s+(1-\epsilon)ds} \left(\frac{m}{dn}\right)^{\epsilon ds - 2s}.$$
(9)

Now, for large enough n, the R.H.S. of (9) vanishes as we have $\frac{m}{dn} = O(\frac{1}{\log n})$

3.2.2 Iterative decoding algorithm

Remember that during the recall phase we are given an n-length observation vector y which is noisy version of one of the message vectors from the dataset \mathcal{M} , i.e.,

$$y = x + e$$
, for some $x \in \mathcal{M}$. (10)

Expander decoding algorithm

```
Input: The vector \mathbf{z} = B\mathbf{e} and the matrix B.

1: Define \mathcal{N}_j := \{i \in [m] : B_{i,j} \neq 0\} \ \forall \ j \in [n].

2: Initialize \widehat{\mathbf{e}} = 0.

3: if \mathbf{z} = B\widehat{\mathbf{e}} then

4: End the decoding and output \widehat{\mathbf{e}}.

5: else

6: Find an index j \in [n] such that the multiset \{\frac{g_i}{B_{i,j}}\}_{i \in \mathcal{N}_j} has at least (1 - 2\epsilon)d identical elements, say \delta. Here, g_i is the gap (cf. (12)) of the constraint defined by the ith row of B.

7: Set \widehat{e}_j \leftarrow \widehat{e}_j + \delta and go to 2.

8: end if
```

Figure 1: Recovery algorithm for sparse vector from expander graphs based measurement matrix [15].

Assuming that we have exactly learnt the $m \times n$ matrix B during the learning phase of the associative memory (as described in Sec. 3.1), we obtain an m-length vector as follows.

$$z = By = B(x + e) = Be, \tag{11}$$

where the last equality follows as we have $x \in \mathcal{M} = \{x \in \mathbb{R}^n : Bx = 0\}$. Note that we have reduced the problem of recovery of the correct message vector x from y to the task of recovering e from z. Assuming that the error vector e satisfies certain sparsity constraint, the latter problem is exactly the problem of recovering the sparse vector e from its linear measurements via the measurement matrix B. As shown in Proposition 1, w.h.p., the matrix B corresponds to the adjacency matrix of an expander graph. In [15], Jafarpour et al. have adapted the iterative error correction algorithm for expander codes from [29] to the problem of sparse recovery problem when the measurement matrix corresponds the adjacency matrix of a good expander graph. Here we propose to employ this iterative algorithm to recover e from e. The algorithm requires calculation of e0 for each of the linear constraints defined by the matrix e1 (or rows of the matrix e3) which we formally define below.

Definition 1. Let e be an error vector and z = Be. Given an estimate \hat{e} for e, for each linear constraint indexed by $i \in [m]$, we define a gap g_i as follows.

$$g_i = z_i - \sum_{j=1}^n B_{i,j} \hat{e}_j. {12}$$

We describe the algorithm in Fig. 1 and present the theoretical guarantees for the performance of the algorithm from [15] as follows.

Proposition 2 ([15]). Let B be an $m \times n$ matrix which is the adjacency matrix for a $(2k, (1-\epsilon)d)$ expander bipartite graph with $\epsilon \leq \frac{1}{4}$. Then, given the measurement vector $\mathbf{z} = B\mathbf{e}$ for any k-sparse vector \mathbf{e} , the expander decoding algorithm (cf. Fig. 1) successfully recovers \mathbf{e} in at most 2k iterations.

We now employ Proposition 2 to characterize the error correction performance of the designed associative memories during the recall phase.

Theorem 5. Let B be the $m \times n$ matrix generated by the sparse-sub-Gaussian model described in Sec. 2.1 and M denote the dataset associated with the matrix B. Then, with probability at least 1 - o(1), the recall phase based on the iterative decoding algorithm described in Fig. 1 can correct at least $\frac{m^2}{2d^2n}$ adversarial errors.

Proof. It follows from Proposition 1 that with probability at least 1 - o(1), the matrix B corresponds to the adjacency matrix of an $\left(\frac{m^2}{d^2n}, (1-\epsilon)d\right)$ -expander graph. Combining the expansion parameters for this expander graph with the result in Proposition 2, we obtain that the iterative decoding algorithm (cf. Fig. 1) can recover the error vector e from z = Be as long as e has at most $\frac{m^2}{2d^2n}$ non-zero coordinates. Given e and e, it is straightforward to obtain the correct message vector as e and e are e. This completes the proof. \Box

4 Proof sketch of Theorem 3: Associative memory storing binary vectors

Since the graph defined by B is still an expander (with edge weights $\{+1, -1\}$), for the recall phase we rely on the same expander decoding algorithm. We just want to guarantee that $|\mathcal{M}| = |\{x \in \{\pm 1\}^n : Bx = 0\}|$ is of size about $\exp(n - \alpha n \log(d \log n)/\log n)$ w.h.p. The algorithm to learn B is same as that of Theorem 1.

Instead of the random model that we have considered in Sec. 2.1.2, consider a random matrix $B \in \{+1,0,-1\}^{m\times n}$ whose each row has independently and uniformly chosen d' nonzero $(\{+1,-1\})$ values. This model allows us to come up with a straight-forward analysis of number of binary vectors in the null-space, while the original model gives the same estimate but with significantly lengthier analysis, that we omit for the interest of space. Note that $d' \sim d\frac{n}{m}$ w.h.p. Now for a randomly and uniformly chosen ± 1 vector \mathbf{y} of length n, and for some constant c' > 0,

$$\mathbb{P}\{B\boldsymbol{y}=0\} = \left(\binom{d'}{\frac{d'}{2}}/2^{d'}\right)^m \ge \left(\frac{1}{c'd'}\right)^{m/2}.$$

This means $\mathbb{E}[|\mathcal{M}|] \ge 2^n \cdot \left(1/(c'd')\right)^{m/2} = 2^{n-\frac{m}{2}\log(c'd')}$. Substituting, $m = c\frac{n}{\log n}$, we get the promised size of \mathcal{M} .

5 Simulation results

Though our main contribution is theoretical, in this section we evaluate the proposed associative memory on synthetic dataset to verify if our methods works. Only a representative figure is presented here (Fig. 2). We consider three sets of system parameters (m, n, d) for the dataset to be stored. For each set of parameters, we first generate an $m \times n$ random matrix B according to the sparse-sub-Gaussian model (cf. Sec. 2.1). Each non-zero entry of the matrix B is drawn uniformly at random from the set $\{\pm 1, \pm 2, \pm 3\}$. We then generate multiple message vectors which belong to the subspace orthogonal to all the rows of the matrix B and provide the learning phase with these vectors. Given these vectors we employ the dictionary learning based approach described in Sec. 3.1.1 to obtain an estimate \widehat{B} for the matrix B. As guaranteed by Theorem 4, in our simulations, \widehat{B} contains all the rows of the original matrix B (however, in a different order). For all three sets of parameters under consideration, we then utilize the estimate \widehat{B} to evaluate the performance of the expander decoding based recall phase (cf. Sec. 3.2). For a fixed number E of errors, we generate 100 error vectors $\mathbf{e} \in \mathbb{R}^n$ with the number of non-zero entries in each error vector equal to E. The non-zero entries in these vectors are uniformly generated from the set $\{\pm 1, \ldots, \pm 4\}$. The positions of the non-zeros entries in each of these vectors are chosen according to a uniform random permutation on the set [n].

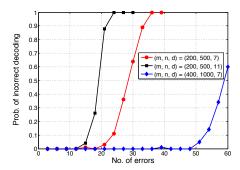


Figure 2: Performance of recall phase for different sets of system parameters.

The performance of the recall algorithm in our simulations is illustrated in Fig. 2 where we plot the fraction of incorrectly recovered error vectors as we increase the number of errors. As expected from Theorem 5, increasing d while keeping m and n fixed degrades the performance of the recall phase. On the other hand, increasing m while keeping d and the ratio $\frac{m}{n}$ fixed improves the performance of the recall phase.

Concluding remarks While we use dictionary learning as a tool in the learning phase, the model of our datasets are subspace models. A large number of datasets on the other hand are also modeled by the *sparse dictionary* model (or union of subspaces). It is of interest to design associative memories, where the datasets are modeled as such. One other possible direction of future research would be to consider a subspace model with a mixture of sparse and dense constraints, which potentially will be inclusive of larger classes of real datasets. For such datasets, under suitable assumption on the generative model, one can potentially employ the techniques of recovering planted sparse vectors in a subspace spanned by dense random sub-Gaussian vectors [5,25] and utilize the recovered sparse constraints to design an iterative recall phase similar to the one presented in this paper. As in the case of [18], the networks (graphs) appearing in our associative memory design share some similarities with the neural networks used for classification tasks. It is an interesting problem to further explore such connections.

References

- [1] R. Adamczak. A note on the sample complexity of the er-spud algorithm by spielman, wang and wright for exact recovery of sparsely used dictionaries. *CoRR*, abs/1601.0204, 2016.
- [2] A. Andoni and P. Indyk. Near-optimal hashing algorithms for approximate nearest neighbor in high dimensions. *Commun. ACM*, 51(1):117–122, Jan. 2008.
- [3] J. Blasiok and J. Nelson. An improved analysis of the er-spud dictionary learning algorithm. *CoRR*, abs/1602.05719, 2016.
- [4] E. J. Candes and T. Tao. Near-optimal signal recovery from random projections: Universal encoding strategies? *IEEE Trans. on Inf. Theory*, 52(12):5406–5425, Dec 2006.
- [5] L. Demanet and P. Hand. Scaling law for recovering the sparsest element in a subspace. *Information and Inference*, 2014.
- [6] David L Donoho. Compressed sensing. IEEE Trans. on Inf. Theory, 52(4):1289–1306, 2006.
- [7] D. Ferro, V. Gripon, and X. Jiang. Nearest neighbour search using binary neural networks. In *Proceedings of IJCNN*, July 2016.
- [8] Vincent Gripon and Claude Berrou. Sparse neural networks with large learning diversity. *IEEE Transactions on Neural Networks*, 22(7):1087–1096, 2011.
- [9] D. J. Gross and M. Mezard. The simplest spin glass. *Nuclear Physics B*, 240(4):431 452, 1984.
- [10] S. Har-Peled, P. Indyk, and R. Motwani. Approximate nearest neighbor: Towards removing the curse of dimensionality. *Theory of computing*, 8(1):321–350, 2012.
- [11] D. O. Hebb. The organization of behavior: A neuropsychological theory. Psychology Press, 2005.
- [12] C. Hillar and N. M. Tran. Robust exponential memory in hopfield networks. *CoRR*, abs/1411.4625, 2014.
- [13] J. J. Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the national academy of sciences*, 79(8):2554–2558, 1982.
- [14] P. Indyk and R. Motwani. Approximate nearest neighbors: Towards removing the curse of dimensionality. In *Proceedings of the Thirtieth Annual ACM Symposium on Theory of Computing*, STOC '98, pages 604–613, New York, NY, USA, 1998. ACM.
- [15] S. Jafarpour, W. Xu, B. Hassibi, and R. Calderbank. Efficient and robust compressed sensing using optimized expander graphs. *IEEE Transactions on Information Theory*, 55(9):4299–4308, Sept 2009.
- [16] S. Jankowski, A. Lozowski, and J. M. Zurada. Complex-valued multistate neural associative memory. *IEEE Transactions on Neural Networks*, 7(6):1491–1496, Nov 1996.
- [17] H. Jégou, M. Douze, and C. Schmid. Product quantization for nearest neighbor search. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(1):117–128, Jan 2011.
- [18] A. Karbasi, A. H. Salavati, and A. Shokrollahi. Convolutional neural associative memories: Massive capacity with noise tolerance. *CoRR*, abs/1407.6513, 2014.

- [19] K. R. Kumar, A. H. Salavati, and A. Shokrollahi. Exponential pattern retrieval capacity with non-binary associative memory. In *2011 IEEE Information Theory Workshop (ITW)*, pages 80–84, Oct 2011.
- [20] A. Mazumdar and A. S. Rawat. Associative memory via a sparse recovery model. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2683–2691. 2015.
- [21] R. J. McEliece and E. C. Posner. The number of stable points of an infinite-range spin glass memory. *Telecommunications and Data Acquisition Progress Report*, 83:209–215, 1985.
- [22] Robert J McEliece, Edward C Posner, Eugene R Rodemich, and Santosh S Venkatesh. The capacity of the hopfield associative memory. *IEEE Transactions on Information Theory*, 33(4):461–482, 1987.
- [23] M. K. Muezzinoglu, C. Guzelis, and J. M. Zurada. A new design method for the complex-valued multistate hopfield associative memory. *IEEE Transactions on Neural Networks*, 14(4):891–899, July 2003.
- [24] Bruno A. Olshausen and David J. Field. Sparse coding with an overcomplete basis set: A strategy employed by v1? *Vision Research*, 37(23):3311 3325, 1997.
- [25] Q. Qu, J. Sun, and J. Wright. Finding a sparse vector in a subspace: linear sparsity using alternating directions. *CoRR*, abs/1412.4659, 2014.
- [26] Tom Richardson and Ruediger Urbanke. *Modern Coding Theory*. Cambridge University Press, New York, NY, USA, 2008.
- [27] A. H. Salavati and A. Karbasi. Multi-level error-resilient neural networks. In 2012 IEEE International Symposium on Information Theory Proceedings (ISIT), pages 1064–1068, July 2012.
- [28] H. Samet. Foundations of Multidimensional and Metric Data Structures (The Morgan Kaufmann Series in Computer Graphics and Geometric Modeling). Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 2005.
- [29] M. Sipser and D. A. Spielman. Expander codes. *IEEE Trans. on Inf. Theory*, 42(6):1710–1722, Nov 1996
- [30] D. A. Spielman, H. Wang, and J. Wright. Exact recovery of sparsely-used dictionaries. In 25th Annual Conference on Learning Theory (COLT), 2012. http://www.columbia.edu/ jw2966/papers/SWW12-pp.pdf.
- [31] F. Tanaka and S. F. Edwards. Analytic theory of the ground state properties of a spin glass. i. ising spin glass. *Journal of Physics F: Metal Physics*, 10(12):2769, 1980.
- [32] J. Wang, H. T. Shen, J. Song, and J. Ji. Hashing for similarity search: A survey. *CoRR*, abs/1408.2927, 2014.
- [33] C. Yu, V. Gripon, X. Jiang, and H. Jégou. Neural associative memories as accelerators for binary vector search. In *Proceedings of Cognitive*, pages 85–89, March 2015.

Appendix

A The modified ER-SpUD algorithm and proof of Theorem 4

In [30], Spielman et al. consider the following problem of exact dictionary learning. Let $D \in \mathbb{R}^{m \times m}$ be an invertible matrix also referred to as the dictionary. Given n observations

$$\mathbf{u}_j = D\mathbf{v}_j \text{ for } j \in [n], \tag{13}$$

the task is to exactly learn the dictionary D and the coefficient matrix

$$V = [\boldsymbol{v}_1, \boldsymbol{v}_2, \dots, \boldsymbol{v}_n] \in \mathbb{R}^{m \times n}.$$

Spielman et al. assume that the coefficient vectors of the observation are randomly generate so that the entries of the coefficient matrix V are independent and identically distributed [30]. In particular, let

$$V_{i,j} = \eta_{i,j} R_{i,j},$$

where $\eta_{i,j} \in \{0,1\}$ and $R_{i,j} \in \mathbb{R}$ are independent random variables. In particular, for some constant α , they assume that

$$\mathbb{P}\left\{\eta_{i,j}=1\right\} = 1 - \mathbb{P}\left\{\eta_{i,j}=0\right\} = \theta \in \left[\frac{2}{m}, \frac{\alpha}{\sqrt{m}}\right],\tag{14}$$

and $R_{i,j}$ is a zero mean sub-Gaussian random variable such that

$$\mathbb{E}[|R_{i,j}|] \ge \frac{1}{10}$$
 and $\mathbb{P}\{|R_{i,j}| \ge t\} \le 2\exp(-t^2/2)$.

This random generative model for the coefficients V is referred to as *Bernoulli-sub-Gaussian model*. Under the Bernoulli-sub-Gaussian model, Spielman et al. show that the dictionary learning problem is well defined. In particular, as long as $n \ge \Omega(m \log m)$, for an alternative representation of the observations

$$U = [\boldsymbol{u}_1, \boldsymbol{u}_2, \dots, \boldsymbol{u}_n] = D'V',$$

where $A' \in \mathbb{R}^{m \times m}$ is an invertible matrix and $V' \in \mathbb{R}^{m \times n}$ is coefficient matrix with the per-column sparsity bounded by that of the original coefficient matrix V, we have

$$D' = D\Pi\Lambda$$

and

$$V' = \Lambda^{-1}\Pi V$$

Here, $\Lambda \in \mathbb{R}^{m \times m}$ and $\Pi \in \mathbb{R}^{m \times m}$ denote a diagonal matrix and a permutation matrix, respectively. This implies that for $n \geq \Omega(m \log m)$, any other representation of the observations which is explained by a square dictionary and the sparsest coefficient vectors have its dictionary and coefficient matrix as some permutation and scaling of the columns and rows of the original dictionary D and the coefficient matrix V, respectively. Furthermore, Spielman et al. also present an algorithm for the exact dictionary learning problem that recovers the dictionary D (up to scaling and permutations of the columns of D) and the coefficient matrix V (up to scaling and permutations of the rows of V) provided that $n = O(m^2 \log^2 m)$ samples. Recently, Adamczak further improve the sample complexity of the dictionary learning algorithm² to $n = O(m \log m)$ [1]. Since we rely on the dictionary learning algorithm in this paper, we describe the algorithm in Fig. 3 and present the exact recovery guarantees from [1].

²In [1], Adamczak analyze a slight modification of the dictionary learning algorithm proposed by Spielman et al.

Modiefied ER-SpUD (DC): Exact recovery of sparsely-used dictionaries using the sum of two columns of U as constraint vectors.

```
Input: n observations U = [\boldsymbol{u}_1, \boldsymbol{u}_2, \dots, \boldsymbol{u}_n] \in \mathbb{R}^{m \times n}.
  1: Initialize the set \mathcal{V} = \emptyset.
  2: for i = 1, ..., n-1 do
           for j = i + 1, ..., n do
  4:
                Let \boldsymbol{r}_{ij} = \boldsymbol{u}_i + \boldsymbol{u}_j.
               Solve minimize \mathbf{w} \in \mathbb{R}^m \| \mathbf{w}^T U \|_1 subject to \mathbf{r}_{ij}^T \mathbf{w} = 1, and set \mathbf{s}_{ij} = \mathbf{w}^T Y \in \mathbb{R}^n.
  5:
                \mathcal{V} = \mathcal{V} \cup \{s_{ij}\}
  6:
           end for
  7:
  8: end for
  9: for i = 1, ..., m do
           Repeat
 10:
           v_i \leftarrow \operatorname{argmin}_{v \in \mathcal{V}} ||v||_0, breaking ties arbitrarily
 11:
 12:
           Until rank([\boldsymbol{v}_1, \boldsymbol{v}_2, \dots, \boldsymbol{v}_i]) = i.
 13:
 14: end for
Output: V = [v_1, v_2, ..., v_m]^T and D = UU^T(VV^T)^{-1}.
```

Figure 3: Description of the dictionary learning algorithm from [1].

Proposition 3. There exists absolute constants $c, \alpha \in (0, \infty)$ such that if

$$\frac{2}{m} \le \theta \le \frac{\alpha}{\sqrt{m}}$$

and V follows the Bernoulli-sub-Gaussian model with parameter θ , then for $n \ge cm \log m$, with probability at least 1 - 1/n the modified ER-SpUD algorithm (cf. Fig. 3) successfully recovers all the rows of V, i.e., multiples of all the rows of U are present among the set V.

A.1 Proof of Theorem 4

In this section we highlight the proof of Theorem 4 which provides the guarantees for the exact recovery of the matrix B using the learning algorithm described in Fig. 3. In [1, Theorem 1.1], Adamczak establishes the analogue of Theorem 4 for $m \times n$ matrices generated by the Bernoulli-sub-Gaussian model (cf. Sec. A). Theorem 4 can be established by suitably modifying the analysis of Adamczak which comprises four main steps as highlighted in [1, Sec. 2.1]. Due to the small differences between the sparse-sub-Gaussian model (cf. Sec. 2.1) for the matrix B considered in this paper and the Bernoulli-sub-Gaussian model from [1], these steps continue to work after small modifications in the analysis. In the rest of this section, we demonstrate this by establishing Lemma 1 for the sparse-sub-Gaussian model which is analogue to [1, Lemma 2.4] for the Bernoulli-sub-Gaussian model. The analogue to other key lemmas from [1] can be similarly obtained.

Let's first define the required notation. In what follows, for $p \ge 1$,

$$\|oldsymbol{v}\|_p := \Big(\sum_{i=1}^m v_i^p\Big)^{1/p}$$

denotes the ℓ_p -norm of the vector $v \in \mathbb{R}^m$. Moreover, we use $B_1^m \subset \mathbb{R}^m$ to denote the set of m-length vectors with unit ℓ_1 -norm, i.e.,

$$B_1^m := \{ \boldsymbol{v} \in \mathbb{R}^m : \| \boldsymbol{v} \|_1 = 1 \}.$$

In [1], Adamczak proves the following concentration result using Bernstein's inequality and Talagrand's contraction principle. Here, we restate this result as it is utilized in the proof of Lemma 1 below.

Proposition 4. [1, Proposition 2.1] Let $R_1, R_2, \ldots, R_n \in \mathbb{R}^m$ and $\xi_1, \xi_2, \ldots, \xi_n \in \{0, 1\}^m$ be two sets of independent random vectors. Assume that for some constant L, we have

$$\mathbb{E}\left[e^{|R_{i,j}|/L}\right] \le 2 \quad \forall \ 1 \le i \le m, \ 1 \le j \le n. \tag{15}$$

Furthermore, assume that we have

$$\mathbb{P}\{\xi_{i,j} = 1\} \le \theta \quad \forall \ 1 \le i \le m, \ 1 \le j \le n. \tag{16}$$

Let $Z_1, Z_2, \dots, Z_n \in \mathbb{R}^m$ be n random vectors defined as follows.

$$Z_{j} = (R_{1,j}\xi_{1,j}, R_{2,j}\xi_{2,j}, \dots, R_{m,j}\xi_{m,j})^{T} \quad \forall \ 1 \le j \le n.$$
(17)

Consider the random variable

$$W = \sup_{\boldsymbol{v} \in B_1^m} \left| \frac{1}{n} \sum_{j=1}^n \left(\boldsymbol{v}^T Z_j - \mathbb{E} \left[\boldsymbol{v}^T Z_j \right] \right) \right|.$$
 (18)

Then, for some universal constant C and every $q \ge \max(2, \log m)$, we have

$$||W||_q \le \frac{C}{n} \left(\sqrt{nq\theta} + q\right) L \tag{19}$$

and

$$\mathbb{P}\left\{W \ge \frac{Ce}{n} \left(\sqrt{nq\theta} + q\right)L\right\} \le e^{-q}.$$
 (20)

Before we proceed, we make the following simple claim about our generative model.

Claim 1. For the random matrix ensemble generated by the sparse-sub-Gaussian model (cf. Sec. 2.1), whenever d = o(m), we have have the following

$$(1 - o(1))\frac{d}{m} \le \mathbb{P}\{\xi_{i,j} = 1\} = 1 - \left(1 - \frac{1}{m}\right)^d \le \frac{d}{m}.$$

We now present the following result which is analogue to [1, Lemma 2.4].

Lemma 1. Let $S \subseteq [n]$ be a fixed subset of size $|S| < \frac{n}{4}$. Let $X \in \mathbb{R}^{m \times n}$ be an $m \times n$ matrix which is generated as follows.

(i) For every $j \in \bar{S} := [n] \setminus S$, we pick d elements uniformly at random from [m] with replacement. Let $\mathcal{N}_j \subseteq [m]$ denote the set of picked elements. Let $R_j = (R_{1,j}, R_{2,j}, \dots, R_{m,j})^T \in \mathbb{R}^m$ be a vector containing i.i.d. sub-Gaussian random variables and $\xi_j \in \{0,1\}^m$ denote the indicator vector for the set $\mathcal{N}_j \subseteq [m]$. Now the j-th column of the matrix X is defined as $X_j = (\xi_{1,j}R_{1,j}, \xi_{2,j}R_{2,j}, \dots, \xi_{m,j}R_{m,j})^T \in \mathbb{R}^m$

(ii) Let $s \leq 2d$. For every $j \in \mathcal{S}$, we pick d elements uniformly at random from the set [m+s] with replacement. Let $\widetilde{\mathcal{N}}_j \subseteq [m+s]$ denote the set of picked elements. We take a subset $\mathcal{N}_j = \widetilde{\mathcal{N}}_j \cap [m]$. Let $R_j = (R_{1,j}, R_{2,j}, \dots, R_{m,j})^T \in \mathbb{R}^m$ be a vector containing i.i.d. sub-Gaussian random variables and $\xi_j \in \{0,1\}^m$ denote the indicator vector for the set $\mathcal{N}_j \subseteq [m]$. Now the j-th column of the matrix X is defined as $X_j = (\xi_{1,j}R_{1,j}, \xi_{2,j}R_{2,j}, \dots, \xi_{m,j}R_{m,j})^T \in \mathbb{R}^m$.

Let X_S denote the sub-matrix of X comprising the columns indexed by the set $S \subseteq [n]$. Then, with probability at least $1 - n^{-8}$, for any $v \in \mathbb{R}^m$, we have

$$\|\boldsymbol{v}^T \boldsymbol{X}\|_1 - 2\|\boldsymbol{v}^T \boldsymbol{X}_{\mathcal{S}}\|_1 \ge \Omega \left(n\mu\sqrt{\frac{\theta}{m}}\right)\|\boldsymbol{v}\|_1,$$
(21)

where $\theta = \frac{d}{m}$ and $\mathbb{E}[|R_{i,j}|] \leq \mu$.

Proof. It follows from Proposition 4 that, with probability at least $1 - n^{-8}$, we have that

$$\sup_{\boldsymbol{v} \in B_1^m} \left| \| \boldsymbol{v}^T X \|_1 - \mathbb{E} [\| \boldsymbol{v}^T X \|]_1 \right| \le C \left(\sqrt{n\theta \log n} + \log n \right)$$

$$\le 2C \sqrt{n\theta \log n}$$

and

$$\sup_{\boldsymbol{v} \in B_1^m} \left| \| \boldsymbol{v}^T X_{\mathcal{S}} \|_1 - \mathbb{E} [\| \boldsymbol{v}^T X_{\mathcal{S}} \|]_1 \right| \le 2C \sqrt{n\theta \log n}.$$

This implies that for any $v \in \mathbb{R}^m$, we have that

$$\|\boldsymbol{v}^T X\|_1 \ge \mathbb{E}[\|\boldsymbol{v}^T X\|]_1 - 2C\sqrt{n\theta \log n}\|\boldsymbol{v}\|_1$$

and

$$\|\boldsymbol{v}^T X_{\mathcal{S}}\|_1 \leq \mathbb{E}[\|\boldsymbol{v}^T X_{\mathcal{S}}\|]_1 + 2C\sqrt{n\theta \log n}\|\boldsymbol{v}\|_1.$$

Combining these two inequalities, we obtain that the following holds for each $v \in \mathbb{R}^m$.

$$\|\boldsymbol{v}^{T}X\|_{1} - 2\|\boldsymbol{v}^{T}X_{\mathcal{S}}\|_{1} \geq \mathbb{E}[\|\boldsymbol{v}^{T}X\|]_{1} - 2\mathbb{E}[\|\boldsymbol{v}^{T}X_{\mathcal{S}}\|]_{1}$$

$$-6C\sqrt{n\theta\log n}\|\boldsymbol{v}\|_{1}$$

$$= \sum_{j\in\mathcal{S}}\mathbb{E}[|\boldsymbol{v}^{T}X_{j}|] + \sum_{j\in\bar{\mathcal{S}}}\mathbb{E}[|\boldsymbol{v}^{T}X_{j}|] - 2\sum_{j\in\mathcal{S}}\mathbb{E}[|\boldsymbol{v}^{T}X_{j}|]$$

$$-6C\sqrt{n\theta\log n}\|\boldsymbol{v}\|_{1}$$

$$= \sum_{j\in\bar{\mathcal{S}}}\mathbb{E}[|\boldsymbol{v}^{T}X_{j}|] - \sum_{j\in\mathcal{S}}\mathbb{E}[|\boldsymbol{v}^{T}X_{j}|]$$

$$-6C\sqrt{n\theta\log n}\|\boldsymbol{v}\|_{1}.$$
(22)

We now bound $\mathbb{E}\left[\left|v^TX_j\right|\right]$ for $j\in\mathcal{S}$. Recall that all the columns of the matrix X indexed by the set \mathcal{S} are identically distributed. Similarly, all the columns of the matrix X indexed by the set $\bar{\mathcal{S}}=[n]\backslash\mathcal{S}$ are identically distributed. In the following, we use $Z=(Z_1,\ldots,Z_m)^T$ and $\hat{Z}=(\hat{Z}_1,\ldots,\hat{Z}_m)^T$ to denote two random vectors with their distribution identical to the columns of the matrix X indexed by the set $\bar{\mathcal{S}}$ and \mathcal{S} , respectively.

$$\mathbb{E}[|\boldsymbol{v}^T\widehat{Z}|] = \mathbb{E}[|\boldsymbol{v}^T(\widehat{Z} + Z - Z)|]$$

$$\leq \mathbb{E}[|\boldsymbol{v}^T Z|] + \mathbb{E}[|\boldsymbol{v}^T (Z - \widehat{Z})|]$$

$$\leq \mathbb{E}[|\boldsymbol{v}^T Z|] + \mu \mathbb{E}[\sum_{i=1}^m |v_i| W_i].$$
 (23)

Here, for $1 \le i \le m$, $W_i = \sum_{l=1}^d Y_i^l$ denotes the sum of d indicator random variable which are defined as follows.

$$Y_i^l = \begin{cases} 1 & \text{with probability } \frac{2d}{m+2d} \frac{1}{m} \\ 0 & \text{with probability } 1 - \frac{2d}{m+2d} \frac{1}{m}. \end{cases}$$
 (24)

Combining (23) and (24), we obtain

$$\mathbb{E}[|\boldsymbol{v}^T \widehat{Z}|] \leq \mathbb{E}[|\boldsymbol{v}^T Z|] + \mu \frac{2d^2}{(m+2d)m} \|\boldsymbol{v}\|_1$$

$$\leq \mathbb{E}[|\boldsymbol{v}^T Z|] + \mu \frac{2d^2}{m^2} \|\boldsymbol{v}\|_1.$$
(25)

Note that we have

$$\mathbb{E}[|\boldsymbol{v}^T X_j|] = \mathbb{E}[|\boldsymbol{v}^T \widehat{Z}|] \quad \forall j \in \mathcal{S}$$
(26)

$$\mathbb{E}[|\boldsymbol{v}^T X_j|] = \mathbb{E}[|\boldsymbol{v}^T Z|] \quad \forall j \in \bar{\mathcal{S}} = [n] \backslash \mathcal{S}.$$
(27)

Therefore, combining (22) and (25), we obtain that for any $v \in \mathbb{R}^m$,

$$\|\boldsymbol{v}^{T}X\|_{1} - 2\|\boldsymbol{v}^{T}X_{\mathcal{S}}\|_{1} \ge (p - 2|\mathcal{S}|)\mathbb{E}[|\boldsymbol{v}^{T}Z|]$$
$$-|\mathcal{S}|\mu \frac{2d^{2}}{m^{2}}\|\boldsymbol{v}\|_{1} - 4C\sqrt{n\theta \log n}\|\boldsymbol{v}\|_{1}.$$
 (28)

Now using $m=c\frac{n}{\log n}, d \leq c'' \log n$ and the lower bound on $\mathbb{E}\big[\big| \boldsymbol{v}^T Z \big|\big]$ from [1, Lemma 2.3], one can argue that

$$\|\boldsymbol{v}^T X\|_1 - 2\|\boldsymbol{v}^T X_{\mathcal{S}}\|_1 \ge \Omega \left(n\mu\sqrt{\frac{\theta}{m}}\right)\|\boldsymbol{v}\|_1.$$
(29)