# On the Role of Channel Capacity in Learning Gaussian Mixture Models

Elad Romanov

ELAD.ROMANOV@GMAIL.COM

The Hebrew University of Jerusalem

BENDORY@TAUEX.TAU.AC.IL

**Tamir Bendory** *Tel Aviv University* 

OR.ORDENTLICH@MAIL.HUJI.AC.IL

Or Ordentlich

The Hebrew University of Jerusalem

Editors: Po-Ling Loh and Maxim Raginsky

### Abstract

This paper studies the sample complexity of learning the k unknown centers of a balanced Gaussian mixture model (GMM) in  $\mathbb{R}^d$  with spherical covariance matrix  $\sigma^2 I$ . In particular, we are interested in the following question: what is the maximal noise level  $\sigma^2$ , for which the sample complexity is essentially the same as when estimating the centers from labeled measurements? To that end, we restrict attention to a Bayesian formulation of the problem, where the centers are uniformly distributed on the sphere  $\sqrt{d}S^{d-1}$ . Our main results characterize the exact noise threshold  $\sigma^2$ below which the GMM learning problem, in the large system limit  $d, k \to \infty$ , is as easy as learning from labeled observations, and above which it is substantially harder. The threshold occurs at  $\frac{\log k}{d} = \frac{1}{2} \log \left(1 + \frac{1}{\sigma^2}\right)$ , which is the capacity of the additive white Gaussian noise (AWGN) channel. Thinking of the set of k centers as a code, this noise threshold can be interpreted as the largest noise level for which the error probability of the code over the AWGN channel is small. Previous works on the GMM learning problem have identified the *minimum distance* between the centers as a key parameter in determining the statistical difficulty of learning the corresponding GMM. While our results are only proved for GMMs whose centers are uniformly distributed over the sphere, they hint that perhaps it is the decoding error probability associated with the center constellation as a channel code that determines the statistical difficulty of learning the corresponding GMM, rather than just the minimum distance.

### 1. Introduction

Gaussian mixture models (GMMs) are widely used in statistics and machine learning. Here, we consider the simplest case of a *spherical*, *balanced* d-dimensional GMM with k-components. Specifically, for centers  $\mathcal{X}_k = (X_1, \dots, X_k) \in \mathbb{R}^{d \times k}$  and variance  $\sigma^2$ , the corresponding GMM, denoted by  $\mathrm{GMM}_{d,k}(\mathcal{X}_k, \sigma^2)$ , is described by the probability distribution  $\mathbf{Y} \sim \mathrm{GMM}_{d,k}(\mathcal{X}_k, \sigma^2)$ :

$$Y = X_{\ell} + \sigma Z$$
,  $\ell \sim \text{Unif}([k])$ ,  $Z \sim \mathcal{N}(0, I)$ , (1)

where  $[k] = \{1, \ldots, k\}$ , and  $\ell \in [k]$  will sometimes be referred to as the *label* of  $\boldsymbol{Y}$  and is statistically independent of  $\boldsymbol{Z}$ . Our focus is on the classical GMM learning problem, where one observes n independent samples  $\boldsymbol{Y}_1, \ldots, \boldsymbol{Y}_n \sim \mathrm{GMM}_{d,k}(\boldsymbol{\mathcal{X}}_k, \sigma^2)$ , and wishes to recover the unknown centers  $\boldsymbol{\mathcal{X}}_k$  (throughout, we always assume that the number of centers k and the variance  $\sigma^2$  are known).

This paper is devoted to studying the fundamental information-theoretic limits of the GMM learning problem, namely, the *sample complexity*: what is the smallest number of samples n one

needs to collect in order to recover the centers (to within some prescribed precision)? The main difficulty in learning the GMM centers is that the samples are unlabeled, and the sample complexity is clearly lower bounded by that of the "genie-aided" setup where each sample is labeled. For sufficiently small noise levels the measurements can be accurately clustered, and the problem is as easy as in the "genie-aided" case, while for large enough noise levels reliable clustering is impossible. The main question we seek to answer here is: what is the critical noise level below which the problem is as statistically easy as in the labeled case, and above which it is significantly harder?

Past works have shown that the *separation* between the centers  $X_1,\ldots,X_k$  has a decisive effect on the statistical difficulty of the problem. Let  $\Delta(\mathcal{X}_k) = \min_{1 \leq i < j \leq k} \|X_i - X_j\|$  be the minimal separation between any two centers. The seminal paper Regev and Vijayaraghavan (2017) has accurately identified the *scaling* of  $\Delta(\mathcal{X}_k)$ , in the large system limit  $k, d \to \infty$ , under which one can estimate the centers (say, to within a small constant precision) using only  $n = \operatorname{poly}(k, d)$  many samples. They show: 1 1) *Upper bound:* If  $\Delta = \Omega(\sigma\sqrt{\log k})$  then the centers may be estimated with  $n = \operatorname{poly}(d, k)$  samples; 2) *Lower bound:* For any  $\gamma(k) = o\left(\sigma\sqrt{\log k}\right)$ , the class of GMMs with minimum separation  $\Delta \geq \sigma\gamma(k)$  is not learnable (in a minimax sense) from  $n = \operatorname{poly}(k, d)$  samples. The upper bound was recently improved by Kwon and Caramanis (2020), who showed that when  $\Delta = \Omega(\sigma\sqrt{\log k})$ , in fact  $n = O(\sigma^2k \cdot \operatorname{polylog}(k))$  samples suffice; this *almost* matches (up to  $\operatorname{polylog}(k)$  factors) the sample complexity for the labeled case. Stated differently, the results above identify the critical noise level scaling for the minimax estimation problem as  $\sigma^2 \sim \frac{\Delta^2(\mathcal{X}_k)}{\log k}$ .

The goal of this paper is to develop a finer grained understanding of the *exact* critical noise level  $\sigma$ , rather than only its scaling. To tackle this ambitious question, we make two modifications with respect to the setup studied in Regev and Vijayaraghavan (2017) and Kwon and Caramanis (2020): 1) Rather than studying the minimax setting with respect to all sets of centers  $\mathcal{X}_k$  with a given  $\Delta(\mathcal{X}_k)$ , we take a Bayesian approach and assume  $\mathcal{X}_k \sim (\mathrm{Unif}(\sqrt{d}\mathcal{S}^{d-1}))^{\otimes k}$ ; 2) We consider a "more forgiving" loss function, which measures the average error in the center reconstruction rather than the maximal error. The rationale behind these modifications will be clarified in the sequel.

Under this setup, we show that the the critical noise level is precisely characterized by the equation  $\frac{1}{2}\log\left(1+\frac{1}{\sigma^2}\right)=\frac{\log k}{d}$ , which is, by no accident, the noise level below which a "typical" constellation  $\mathcal{X}_k$  constitutes a good error correcting code for the AWGN( $\sigma^2$ ) channel (additive white Gaussian noise, with noise variance  $\sigma^2$ ). Our analysis relies *explicitly* on the decodability properties of  $\mathcal{X}_k$ , when thought of as a channel code. This is a "global" property of the constellation, compared to the minimum distance (note that it is well-known that at high coding rate, the minimum distance of a code is not entirely predictive of its error probability, see e.g. Barg and Forney (2002)). Regarding the minimum separation, we remark that, as is to be expected, our results are consistent with Regev and Vijayaraghavan (2017) regarding the required scaling of  $\Delta(\mathcal{X}_k)$  for statistically-efficient learning. Classical results on sphere packing, e.g., Kabatiansky and Levenshtein (1978), imply that if  $\log k/d$  is finite, "typical" constellations under  $\mathcal{X}_k \sim (\mathrm{Unif}(\sqrt{d}\mathcal{S}^{d-1}))^{\otimes k}$  have minimal separation  $\Delta(\mathcal{X}_k) = \Theta(\sqrt{d})$ . Thus, 1) When  $\log k/d = \Theta(1)$  the critical noise level is at  $\sigma^2 = \Theta(1)$ , so in terms of minimal separation,  $\Delta(\mathcal{X}_k)/\sigma = \Theta(\sqrt{\log k})$ ; 2) On the other hand, when  $\log k/d = o(1)$ , the critical noise level is  $\sigma^2 = \Theta(d/\log k)$  and so  $\Delta(\mathcal{X}_k)/\sigma = \Theta(\sqrt{\log k})$ .

Finally, our results *hint* at the possibility of a deeper connection between channel coding and statistical inference: the decodability properties of the set of centers  $\mathcal{X}_k$  (as a channel code) may

<sup>1.</sup> We restrict our attention in this discussion, and throughout the paper, exclusively to an asymptotic regime where  $d,k\to\infty$  together with  $\limsup_{d,k\to\infty}\frac{\log k}{d}<\infty$ .

determine, to an extent, the statistical difficulty of learning the corresponding GMM. The present paper takes a modest first step towards showing such a connection, establishing it for the special case of spherical random codes, whose typical instances posses strong symmetry properties.

### 1.1. Formal Problem Formulation

As mentioned before, we study the large system behavior of the sample complexity under a uniform spherical prior on the centers. Denote the (random) centers by

$$\boldsymbol{\mathcal{X}}_k = (\boldsymbol{X}_1, \dots, \boldsymbol{X}_k) \sim \left( \text{Unif}(\sqrt{d}\boldsymbol{\mathcal{S}}^{d-1}) \right)^{\otimes k}.$$
 (2)

Note that we scale the problem so that  $||X_i|| = \sqrt{d}$  for all  $i \in [k]$ . We observe n measurements,  $Y_1, \ldots, Y_n$ , sampled from the GMM distribution whose centers are  $\mathcal{X}_k$ :

$$\begin{bmatrix} \boldsymbol{Y}_1, \dots, \boldsymbol{Y}_n \mid \boldsymbol{\mathcal{X}}_k \end{bmatrix} \overset{i.i.d.}{\sim} \operatorname{GMM}_{d,k}(\boldsymbol{\mathcal{X}}_k, \sigma^2),$$
 (3)

see also (1). Per standard terminology in signal processing,  $1/\sigma^2$  may be interpreted as the "signal-to-noise ratio" (SNR) per coordinate. Suppose that  $\hat{\boldsymbol{\mathcal{X}}}_k = (\hat{\boldsymbol{X}}_1, \dots, \hat{\boldsymbol{X}}_k)$  is an estimator of  $\boldsymbol{\mathcal{X}}_k$ , computed from the measurements. The model admits the following Markov chain structure:

$$\mathcal{X}_k = (X_1, \dots, X_k) \longrightarrow (Y_1, \dots, Y_n) \longrightarrow \hat{\mathcal{X}}_k = (\hat{X}_1, \dots, \hat{X}_k).$$
 (4)

At this point it is instructive to think about the much simpler estimation problem, where each measurement  $Y_i$  is observed with its label  $\ell_i \in [k]$ , and every center is observed exactly n/k times. For this problem, the optimal mean squared error (MSE) in the reconstruction of each center is  $d^{-1}\mathbb{E}||X_i - \hat{X}_i||^2 = k\sigma^2/n$  (to leading order in k/n), and is attained for example, by the sample mean. In the GMM estimation problem the samples are not labeled, and furthermore, the number of times each center appears in the measurements is a  $\operatorname{Binomial}(n,1/k)$  random variable. While the mean of this random variable is indeed n/k, some centers will appear fewer times. In particular, when  $n = o(k \log k)$  some of the centers are likely to not appear even once (coupon collecting). To circumvent the issues arising due to this effect, and focus our study on the problem of dealing with the lack of labels, we measure the discrepancy between  $\mathcal{X}_k$  and  $\hat{\mathcal{X}}_k$ , by the loss function

$$\mathcal{L}_{\text{avg}}(\boldsymbol{\mathcal{X}}_k, \hat{\boldsymbol{\mathcal{X}}}_k) = \frac{1}{k} \sum_{i=1}^k d^{-1} \text{dist}^2(\boldsymbol{X}_i, \hat{\boldsymbol{\mathcal{X}}}_k) := \frac{1}{k} \sum_{i=1}^k \min_{1 \le j \le k} d^{-1} \|\boldsymbol{X}_i - \hat{\boldsymbol{X}}_j\|^2.$$
 (5)

In words: the average normalized squared distance between a center  $X_i$  and the list  $\hat{X}_k$ . As we shall see, under this loss function it is possible to obtain a risk of  $k\sigma^2/n$  for  $\sigma$  below the critical noise level and n/k large enough. In contrast, the more restrictive max-loss function  $\mathcal{L}_{\max}(\mathcal{X}_k, \hat{\mathcal{X}}_k) = \max_{1 \leq i \leq k} d^{-1} \mathrm{dist}^2(X_i, \hat{\mathcal{X}}_k)$  considered in much of the prior work, does not decay with n in the regime  $n = o(k \log k)$ , regardless of the noise level, due to the non-uniform empirical distribution of the center indices. Under  $\mathcal{L}_{\mathrm{avg}}$ , on the other hand, to achieve  $\varepsilon$  error, it suffices to estimate only a fraction  $1 - O(\varepsilon)$  of the centers within error  $O(\varepsilon)$ , having the remaining centers incur an error O(1). Thus, the effect of non-uniform label empirical distribution is bypassed by this loss function. Under our formulation of the GMM learning problem, the goal is to construct an estimation

rule  $\hat{\mathcal{X}}_k$ :  $(\mathbb{R}^d)^n \to \mathbb{R}^{d \times k}$  ("algorithm") so to minimize the risk:  $\mathbb{E}\mathcal{L}_{avg}(\mathcal{X}_k, \hat{\mathcal{X}}_k(Y_1, \dots, Y_n))$ .

Importantly, the expectation is taken over the randomness in both the sample generating process given the centers (3), as well as the center prior distribution (2), whose joint distribution adheres to the Markov chain structure in (4). We study the information-theoretic limits of the aforementioned problem. Consider the minimum attainable risk over all estimation laws  $\mathcal{X}_k$ :

$$R_n = \inf_{\hat{\boldsymbol{\mathcal{X}}}_k} \mathbb{E} \mathcal{L}_{avg}(\boldsymbol{\mathcal{X}}_k, \hat{\boldsymbol{\mathcal{X}}}_k(\boldsymbol{Y}_1, \dots, \boldsymbol{Y}_n)).$$
 (6)

For a fixed precision level  $\varepsilon > 0$ , define the *sample complexity*,

$$n_{\varepsilon}^* = n_{\varepsilon}^*(d, k, \sigma^2) = \min\{n : R_n \le \varepsilon\}.$$
 (7)

Importantly, (6) and (7) make no assumptions about the computational difficulty of implementing  $\hat{\mathcal{X}}_k: (\mathbb{R}^d)^n \to \mathbb{R}^{d \times k}$ , and in particular are not restricted to computational efficient algorithms (poly(d, k)) runtime). Throughout, computational considerations shall be completely neglected.

### 1.2. Main Results

As our analysis relies on viewing the centers as a code for the AWGN channel, the problem's rate

$$R_{d,k} := \frac{\log k}{d} \,, \tag{8}$$

and the decreasing function  $C:(0,\infty)\to(0,\infty)$ 

$$C(\sigma^2) = \frac{1}{2} \log \left( 1 + \frac{1}{\sigma^2} \right) , \tag{9}$$

characterizing the AWGN( $\sigma^2$ ) channel capacity, will play a key role. Throughout the paper, we couple the noise level  $\sigma^2$  to k and d by the parameter  $\beta \in (0, \infty)$  via the equation

$$\mathsf{R}_{d,k} = \mathsf{C}(\beta\sigma^2) \,. \tag{10}$$

When  $\beta > 1$ , the rate is smaller than the capacity; when  $\beta < 1$ , it is larger. This parametrization will turn out particularly useful in the statement of the results and their derivations.

We restrict attention to the large-system limit, where  $d, k \to \infty$ , and denote the *limiting rate* by

$$\mathsf{R} = \lim_{d \to \infty} \mathsf{R}_{d,k} \in [0, \infty) \,. \tag{11}$$

We distinguish between two asymptotic regimes:

- (Positive Rate, R > 0):  $\sigma^2 \in (0, \infty)$  is a fixed constant. In particular,  $k = e^{\Theta(d)}$ . (Zero Rate, R = 0):  $\sigma^2 \to \infty$ . So that also  $k \to \infty$ , we also require impose  $\sigma^2 = o(d)$ .

We remark that, since we are interested in estimation to *finite precision*  $\varepsilon$  in Theorems 1 and 2 below, the asymptotic regime  $\log k = \omega(d)$ , namely when the number of centers k is super-exponential in d, becomes rather uninteresting. Indeed, for super-exponential k one may simply take  $\hat{\mathcal{X}}_k$  to be some fixed  $\sqrt{\varepsilon d}$ -net of the sphere  $\sqrt{d}\mathcal{S}^{d-1}$ , which can be of size  $(O(1/\varepsilon))^{d/2} \ll k$ . Clearly,  $\mathcal{L}_{avg}(\mathcal{X}_k, \hat{\mathcal{X}}_k) \leq \varepsilon$  for any  $\mathcal{X}_k$ , so under such asymptotics  $n_{\varepsilon}^* = 0$  exactly.

Our first main result states that when the rate is below the channel capacity,  $\mathcal{X}_k$  is learnable at essentially the same sample complexity as in the labeled case.

**Theorem 1** *Suppose that*  $\beta > 1$ *. Then* 

$$e^{-2\mathsf{R}} \le \lim_{\varepsilon \to 0} \lim_{d \to \infty} \frac{n_{\varepsilon}^*}{(\sigma^2 k/\varepsilon)} \le 1$$
 (12)

Theorem 1 implies that when the rate is below the channel capacity, for every fixed and small precision  $\varepsilon>0$ , and for d large, the sample complexity scales like  $n=C\sigma^2k/\varepsilon$ , where  $C\in[e^{-2R},1]$ . Remarkably, when k is sub-exponential in d (R = 0) the pre-factor C is precisely 1. Thus, the sample complexity of the GMM learning problem is exactly the same as that of the labeled case, up to lower order terms in  $1/\varepsilon$ , and asymptotically  $(d\to\infty)$  vanishing correction terms.

Our second main result states that above the capacity, the sample complexity is super-linear:

**Theorem 2** Suppose that  $\beta < 1$ . Then for any fixed sufficiently small  $\varepsilon < \varepsilon_0(R)$ ,

$$\lim_{d \to \infty} \frac{n_{\varepsilon}^*}{\sigma^2 k/\varepsilon} = \infty. \tag{13}$$

*Moreover, the following quantitative bounds hold for all sufficiently small*  $\varepsilon < \varepsilon_0(R)$ :

1. If R > 0 then

$$\frac{n^*}{\sigma^2 k} = \Omega_{\varepsilon,\beta,\mathsf{R}} \left( \sqrt{\frac{\log k}{\log \log k}} \right) \,. \tag{14}$$

2. If R = 0 then

$$\frac{n^*}{\sigma^2 k} = \Omega_{\varepsilon,\beta} \left( \min \left\{ \sqrt{\frac{\log k}{\log \log k}}, \sqrt{\frac{d}{\log k}} \right\} \right). \tag{15}$$

Theorems 1 and 2 together reveal a dichotomy: *precisely* at the channel capacity ( $\beta = 1$ ), the large-system behavior of the sample complexity undergoes a *phase-transition*, from a linear growth in  $\sigma^2 k$ , as in the labeled case, to super-linear growth.

## 1.3. Prior Art

The problem of estimating the parameters of a Gaussian mixture model has a long and rich history, going back to the pioneering work of Pearson (1894). We briefly mention some pointers to the literature, though we emphasize that the list below is *not exhaustive by any means*.

The first work to highlight the importance of minimum separation in learning GMMs is Dasgupta (1999), who gave a poly-time algorithm assuming (in the spherical balanced case)  $\Delta = \Omega(\sigma\sqrt{d})$ . Subsequent works have gradually improved upon the required bound on  $\Delta$ . Early incarnations include Arora and Kannan (2001); Vempala and Wang (2004); Achlioptas and McSherry (2005); Dasgupta and Schulman (2007); Kannan et al. (2008), which culminated in a bound  $\Delta = \Omega(\sigma k^{1/4})$  as sufficient for estimation in polynomial time. This barrier was broken only fairly recently Diakonikolas et al. (2018); Hopkins and Li (2018); Kothari et al. (2018), who have shown that separation  $\Delta = \Omega(\sigma k^{\gamma})$  suffices for polynomial-time learnability, for *any* constant  $\gamma > 0$ .

As for statistical lower bounds, it is known that in the absence of a separation condition,  $n = \exp(k)$  samples are generally necessary to learn the parameters of a GMM Moitra and Valiant (2010); Hardt and Price (2015). The work Regev and Vijayaraghavan (2017) has shown that

separation  $\Delta = \Omega(\sigma\sqrt{\log k})$  is a *sufficient and necessary* condition for learning GMMs with  $n = \operatorname{poly}(k,d)$  samples; the algorithm they proposed to prove their upper bound has exponential runtime. Kwon and Caramanis (2020) have recently improved their upper bound on the sample complexity, and have show that in fact  $n = O(\sigma^2 k \cdot \operatorname{polylog}(k))$  samples suffice, which *almost* matches the trivial lower bound of  $n = \Omega(\sigma^2 k)$ . Their analysis consists of two components: 1) An exponential-time initialization scheme, that finds points *sufficiently close* to the true centers, based on the results of Ashtiani et al. (2018); 2) New local convergence and finite-sample guarantees for (a slightly modified version of) the well-known Expectation Maximization (EM) algorithm. To our knowledge, the problem of learning  $\Delta = \Omega(\sigma\sqrt{\log k})$ -separated GMMs in polynomial time, or proving that this cannot be done (the existence of a computational-statistical gap) is still open.

Another line of work circumvents the minimal separation requirement, by instead restricting attention to "typical" problem instances, an approach much in line with the results of the present paper, and in the context of learning GMMs dates, to the best of our knowledge, to the study Srebro et al. (2006). In the papers Hsu and Kakade (2013); Bhaskara et al. (2014); Goyal et al. (2014); Anderson et al. (2014); Anandkumar et al. (2014); Ge et al. (2015), it is shown that when the center configuration satisfies certain algebraic *non-degeneracy* conditions, methods based on tensor decomposition may be used to recover the centers; such non-degenerate configurations are highly abundant when d is large relative to k, specifically  $k \le d^{O(1)}$ .

Lastly, a different line of work considers learning GMMs by means of density estimation, that is, given samples  $Y_1, \ldots, Y_n$  one has to construct a density f which is close to  $\mathrm{GMM}_{d,k}(\mathcal{X}_k, \sigma^2)$  in, e.g., total variation distance. This problem may be considered in either in the setting of proper learning (f has to be a k-component GMM) or improper learning (no such restriction), see for example Feldman et al. (2006); Kalai et al. (2010); Chan et al. (2014); Suresh et al. (2014); Li and Schmidt (2017); Diakonikolas et al. (2019); Ashtiani et al. (2018). For well-seaprated spherical GMMs,  $\Delta = \Omega(\sigma\sqrt{\log k})$ , guarantees for proper distribution estimation may be translated to error bounds on the centers, see Regev and Vijayaraghavan (2017); Kwon and Caramanis (2020).

Our proof program closely follows that of Romanov et al. (2021), which studied the sample complexity of the multi-reference alignment (MRA) problem. MRA is a *particular* instance of a GMM, with exactly k=d components corresponding to different shifted versions of the same signal. While, similarly to Romanov et al. (2021), the proof of our lower bound uses the mutual information method Polyanskiy and Wu (2014), here the mutual information is upper bounded using the I-MMSE relation rather than the Fano-based argument of Romanov et al. (2021). More importantly, the proof of the upper bound here requires overcoming several significant hurdles not present in the MRA model. In particular, while in MRA we always have k=d, in the GMM problem k may be much greater, and even exponential in the dimension. Furthermore, in MRA there is a single signal to be estimated and all measurement are informative for its estimation. Here, on the other hand, many centers must be estimated, which significantly complicates the first step of our reconstruction algorithm with respect to that used in Romanov et al. (2021).

**Paper outline.** In Section 2 we provide brief background on channel coding and random spherical codes, which shall be used in the analysis to follow. In Section 3 we outline the proof of the lower bound in Theorems 1 and 2. In Section 4 we outline the proof of the upper bound in Theorem 1. To keep within the space constraint, most of the technical details are deferred to the Appendix.

## 2. Background on Channel Coding

A key message of this paper is the following: the centers  $\mathcal{X}_k$  are learnable at linear sample complexity exactly in the regime where the constellation  $\mathcal{X}_k = (X_1, \dots, X_k)$  defines (with high probability) a good codebook for the AWGN channel with noise variance  $\sigma^2$ . Throughout the analysis, the connection to the decoding capabilities of  $\mathcal{X}_k$  will be instrumental. In this section, we briefly survey the required background from information and coding theory. We refer the reader to Cover and Thomas (2012), Gallager (1968) and Polyanskiy and Wu (2014) for a comprehensive treatment.

A coding scheme for sending  $\log k$  nats over the d-dimensional AWGN channel consists of a codebook and a decoder. The codebook is a set of k codewords  $\mathcal{C} = (C_1, \ldots, C_k) \in \mathbb{R}^{d \times k}$ , where codeword  $C_i$  encodes message  $1 \leq i \leq k$ , and all codewords satisfy  $\|C_i\|^2 \leq d$ . The code's rate is  $\mathsf{R}_{d,k} = \frac{\log k}{d}$ . The decoder  $\mathsf{Dec} : \mathbb{R}^d \to [k]$  is a mapping from channel outputs to messages. It is often convenient to allow the decoder to output symbols in  $[k] \cup \{\#\}$ , where the special symbol # corresponds to a declared decoding error.

The decoding error associated with message  $i \in [k]$ , for a given a codebook-decoder pair, is

$$P_{e,i}(\sigma^2|\mathcal{C}, \mathsf{Dec}) = \Pr\left(i \neq \mathsf{Dec}(X_i + \sigma Z)\right),$$
 (16)

and the average error over all messages is

$$P_{e,avg}(\sigma^2|\mathcal{C}, \mathsf{Dec}) := \frac{1}{k} \sum_{i=1}^k P_{e,i}(\sigma^2|\mathcal{C}, \mathsf{Dec}) = \Pr_{\ell \sim \mathsf{Unif}([k])} (\ell \neq \mathsf{Dec}(X_\ell + \sigma Z)) \ . \tag{17}$$

For a given codebook C, the optimal decoder, in the sense of smallest average error, is clearly given by the maximum a posteriori probability (MAP) rule

$$DecOpt(\mathbf{Y}) = \underset{i \in [k]}{\operatorname{argmax}} \Pr(\ell = i \mid \mathbf{Y}, \mathbf{C}) = \underset{1 \le i \le k}{\operatorname{argmin}} \|\mathbf{Y} - \mathbf{C}_i\|^2,$$
(18)

where ties are broken arbitrarily. Accordingly, we define the error of the codebook C, and the corresponding individual errors as

$$\rho_{\text{avg}}(\sigma^2|\mathcal{C}) = P_{e,avg}(\sigma^2|\mathcal{C}, \text{DecOpt}), \quad \rho_i(\sigma^2|\mathcal{C}) = P_{e,i}(\sigma^2|\mathcal{C}, \text{DecOpt}).$$
 (19)

In communication theory, one is interested in designing coding schemes with large rate and small error probability. We say a rate  $\mathsf{R} \in (0,\infty)$  is *achievable* if there exists a sequence  $(d \to \infty)$  of codebooks  $\mathcal{C} \in \mathbb{R}^{d \times k}$  such that  $\lim_{d \to \infty} \mathsf{R}_{d,k} = \mathsf{R}$  and  $\lim_{d \to \infty} \rho(\sigma^2 | \mathcal{C}) = 0$ . Shannon's celebrated channel coding theorem gives a precise characterization of all the achievable rates:

**Theorem 3 (Channel coding theorem, AWGN channel)** Fix  $\sigma^2$ , and let  $C(\cdot)$  be given in (9).

- 1. (Achievability). Any rate  $R < C(\sigma^2)$  is achievable.
- 2. (Converse). No rate  $R > C(\sigma^2)$  is achievable.

The achievability part of the channel coding theorem is typically proved using a random coding argument with respect to the ensemble of i.i.d. Gaussian codebooks. However, it can also be proved using the ensemble of random spherical codebooks,  $\mathcal{C} = \mathcal{X}_k = (X_1, \dots, X_k) \sim \mathrm{Unif}(\sqrt{d}\mathcal{S}^{d-1})^{\otimes k}$ . In fact, the latter ensemble results in a favorable decay of the error probability with d, Shannon (1959). We denote the decoding error, averaged over the codebook ensemble, by

$$\rho_{\text{avg}}(\sigma^2) = \mathbb{E}[\rho_{\text{avg}}(\sigma^2 | \boldsymbol{\mathcal{X}}_k)] \stackrel{(\star)}{=} \mathbb{E}[\rho_i(\sigma^2 | \boldsymbol{\mathcal{X}}_k)], \qquad (20)$$

where  $(\star)$  holds since each  $X_i$  has the same distribution.

**Proposition 4** Let  $\beta > 1$  be fixed. Suppose that  $d, k \to \infty$ , with  $R_{d,k} = C(\beta \sigma^2)$ , so that either: 1)  $\sigma^2$  fixed; or 2)  $\omega(1) = \sigma^2 = o(d)$ . Then  $\lim_{d\to\infty} \rho_{avg}(\sigma^2) = 0$ .

While Proposition 4 is well-known when  $\sigma^2$  is fixed (positive rate) Shannon (1959), the case of  $\omega(1) = \sigma^2 = o(d)$  has not been mainstreamed. We provide a self-contained proof of Proposition 4 in Appendix, Section A, since it will serve as the baseline for the derivations that follow.

### 3. Proof of Lower Bounds

Our proof of the lower bounds in Theorems 1 and 2 uses a standard framework for proving estimation lower bounds (e.g., (Polyanskiy and Wu, 2014, Chapter 28)).

Suppose  $\hat{\mathcal{X}}_k = \hat{\mathcal{X}}_k(Y_1, \dots, Y_n)$  attains  $\mathbb{E}\mathcal{L}_{avg}(\mathcal{X}_k, \hat{\mathcal{X}}_k) \leq \varepsilon$ . Consider the Markov chain (4). By the data processing inequality (DPI) (Polyanskiy and Wu, 2014, Theorem 2.5),

$$I(\boldsymbol{\mathcal{X}}_k; \hat{\boldsymbol{\mathcal{X}}}_k) \le I(\boldsymbol{\mathcal{X}}_k; \boldsymbol{Y}_1, \dots, \boldsymbol{Y}_n).$$
 (21)

We lower bound the LHS of (21) in terms of  $\varepsilon$  and upper bound the RHS in terms of n and  $\sigma^2$ . Starting with  $I(\mathcal{X}_k; \hat{\mathcal{X}}_k)$ , clearly,

$$I(\boldsymbol{\mathcal{X}}_{k}; \hat{\boldsymbol{\mathcal{X}}}_{k}) \ge \min_{P_{\boldsymbol{\mathcal{D}}|\boldsymbol{\mathcal{X}}_{k}} : \mathbb{E}\mathcal{L}_{avg}(\boldsymbol{\mathcal{X}}_{k}, \boldsymbol{\mathcal{D}}) \le \varepsilon} I(\boldsymbol{\mathcal{X}}_{k}; \boldsymbol{\mathcal{D}}),$$
(22)

where we minimize the mutual information (MI) over all conditional distributions of random variables  $\mathcal{D} = (D_1, \dots, D_k) \in \mathbb{R}^{d \times k}$ , under the expected loss constraint  $\mathbb{E}\mathcal{L}_{avg}(\mathcal{X}_k, \mathcal{D}) \leq \varepsilon$ . The minimization (22) is an instance of a *rate-distortion* problem, typically encountered when studying the information-theoretic limits of lossy compression (Polyanskiy and Wu, 2014, Chapter 25).

One complication that arises when attempting to solve the optimization problem in (22) is that the distortion measure,  $\mathcal{L}_{\text{avg}}(\boldsymbol{\mathcal{X}}_k, \boldsymbol{\mathcal{D}}) = \frac{1}{dk} \sum_{i=1}^k \min_{1 \leq j \leq k} \|\boldsymbol{X}_i - \boldsymbol{D}_j\|^2$  is somewhat nonstandard. If instead we had used the quadratic loss,  $\frac{1}{dk} \|\boldsymbol{\mathcal{X}}_k - \boldsymbol{\mathcal{D}}\|_F^2 = \frac{1}{dk} \sum_{i=1}^k \|\boldsymbol{X}_i - \boldsymbol{D}_i\|^2$ , the resulting optimization problem would essentially lend itself to the classical problem of computing the Gaussian quadratic rate-distortion function (RDF), which admits the solution  $\frac{dk}{2} \log(1/\varepsilon)$ .

The loss  $\mathcal{L}_{avg}$  differs from the standard quadratic loss in that it allows for k additional degrees of freedom: every  $i \in [k]$  is matched to the best index  $j_i = \operatorname{argmin}_{j \in [k]} \| \boldsymbol{X}_i - \boldsymbol{D}_{j_i} \|$ . Since the entropy of the k-tuple  $(j_1, \ldots, j_k)$  is at most  $k \log k$  nats, the RDF for  $\mathcal{L}_{avg}$  must be at most  $k \log k$  nats away from the RDF for the standard quadratic loss. We prove in Appendix, Section B.1:

**Lemma 5** Consider the Markov chain (4), with  $\mathbb{E}\mathcal{L}_{avg}(\boldsymbol{\mathcal{X}}_k, \hat{\boldsymbol{\mathcal{X}}}_k) \leq \varepsilon$ . For universal  $c_0 > 0$ ,

$$I(\boldsymbol{\mathcal{X}}_k; \hat{\boldsymbol{\mathcal{X}}}_k) \ge \frac{dk}{2} \log(1/\varepsilon) - dk \log\left(1 + c_0(\varepsilon d)^{-1/2}\right) - k \log k.$$

Next, we upper bound  $I(\mathcal{X}_k; Y_1, \dots, Y_n)$ , starting with a trivial bound. Let  $\ell = (\ell_1, \dots, \ell_n)$  be the random labels, such that  $Y_j = X_{\ell_j} + \sigma Z_j$ . By the DPI,  $I(\mathcal{X}_k; Y_1, \dots, Y_n) \leq I(\mathcal{X}_k; Y_1, \dots, Y_n, \ell)$ . Now, given  $\ell$ , the mapping  $\mathcal{X}_k \mapsto (Y_1, \dots, Y_n)$  simply corresponds to k parallel Gaussian channels, each used on average n/k times. Thus, as we formally prove in Appendix, Section B.3,

**Lemma 6** The following holds:

$$I(\boldsymbol{\mathcal{X}}_k; \boldsymbol{Y}_1, \dots, \boldsymbol{Y}_n) \le I(\boldsymbol{\mathcal{X}}_k; \boldsymbol{Y}_1, \dots, \boldsymbol{Y}_n, \boldsymbol{\ell}) \le \frac{dk}{2} \log \left(1 + \frac{n}{k\sigma^2}\right).$$
 (23)

Consequently, combining with (21) and Lemma 5,

$$\lim_{d\to\infty} \frac{n_{\varepsilon}^*}{k\sigma^2} \ge e^{-2\mathsf{R}} \varepsilon^{-1} - 1.$$

The bound (23) misses a *crucial* aspect of our problem: the observations are *not* labeled. We next derive a bound which does capture this effect, though at the loss of the "correct" dependence on n.

Observe that  $Y_1, \ldots, Y_n$  are conditionally independent given  $\mathcal{X}_k$ . That is: the "channel" mapping the set of centers to samples is *memoryless*. It is an elementary fact (Polyanskiy and Wu, 2014, Theorem 5.1) that in this case, the MI is subadditive

$$I(\boldsymbol{\mathcal{X}}_k; \boldsymbol{Y}_1, \dots, \boldsymbol{Y}_n) \le \sum_{i=1}^n I(\boldsymbol{\mathcal{X}}_k; \boldsymbol{Y}_i) = n \cdot I(\boldsymbol{\mathcal{X}}_k; \boldsymbol{Y}).$$
 (24)

While this bound fails to correctly capture the dependence of  $I(\mathcal{X}_k; Y_1, \dots, Y_n)$  on n when  $n/(k\sigma^2)$  is large, it does suffice for establishing the phase transition of the sample complexity that we seek here. We proceed to bounding the single-sample MI,  $I(\mathcal{X}_k; Y)$ , a much more manageable object. Let  $\ell \sim \mathrm{Unif}([k])$  be the random label of Y. Using the MI chain rule both ways,

$$I(\mathcal{X}_k, \ell; \mathbf{Y}) = I(\mathcal{X}_k; \mathbf{Y}) + I(\ell; \mathbf{Y}|\mathcal{X}_k) = I(\ell; \mathbf{Y}) + I(\mathcal{X}_k; \mathbf{Y}|\ell).$$

Now,  $I(\ell, \boldsymbol{Y}) = 0$  (since  $\{\boldsymbol{X}_i\}_{i=1}^k$  are identically distributed, so  $\boldsymbol{Y}$  does not depend on  $\ell$ ). Similarly,  $I(\ell; \boldsymbol{Y} | \boldsymbol{\mathcal{X}}_k) = H(\ell | \boldsymbol{\mathcal{X}}_k) - H(\ell | \boldsymbol{\mathcal{X}}_k, \boldsymbol{Y})$ , and  $H(\ell | \boldsymbol{\mathcal{X}}_k) = H(\ell) = \log k$  Furthermore,  $I(\boldsymbol{\mathcal{X}}_k; \boldsymbol{Y} | \ell) = I(\boldsymbol{X}_\ell; \boldsymbol{Y} | \ell) \leq \mathsf{C}(\sigma^2) d$ , as the AWGN channel capacity  $\mathsf{C}(\sigma^2)$  upper bounds  $I(\boldsymbol{X}; \boldsymbol{X} + \sigma \boldsymbol{Z})/d$  for any random variable on  $\mathbb{R}^d$  with second moment  $\mathbb{E} \|\boldsymbol{X}\|^2 \leq d$ . Combining these equalities and estimates and rearranging, we obtain

$$I(\boldsymbol{\mathcal{X}}_k; \boldsymbol{Y}) \le \mathsf{C}(\sigma^2)d - \log k + H(\ell|\boldsymbol{\mathcal{X}}_k, \boldsymbol{Y}). \tag{25}$$

In light of (25), it remains to estimate  $H(\ell|\mathcal{X}_k, \mathbf{Y})$ , to be interpreted as the remaining uncertainty in a message  $\ell$  that is sent across the channel, given the output  $\mathbf{Y}$  as well as the known codebook  $\mathcal{X}_k$ . To that end, consider the non-increasing mapping  $\beta \mapsto \varphi(\beta) = H(\ell|\mathcal{X}_k, \mathbf{Y})$  (recall that larger  $\beta$  corresponds to smaller  $\sigma$ ). Since a typical realization of  $\mathcal{X}_k$  results in a code whose error vanishes when  $\beta > 1$ , Fano's inequality implies that  $\varphi(\beta)|_{\beta>1} = o(\log(k))$ . Thus, for  $\beta < 1$ , we have that  $\varphi(\beta) = -\int_{\beta}^{1+\delta} \varphi'(s) ds + o(\log(k))$ , for any  $\delta > 0$ . Using the I-MMSE formula Guo et al. (2005), a remarkable connection between information and estimation under Gaussian channels, the derivative  $\varphi'(\beta)$  can be expressed as the minimum MSE (MMSE) in estimating  $\mathbf{X}_\ell$  from  $\mathbf{Y}$ . Finally, we upper bound the MMSE by the optimal MSE for linear estimation, resulting in the following lemma, whose full proof appears in Appendix, Section B.4. We denote by  $\mathsf{C}^{-1}$  the inverse of (9), and by  $h_b(p) = p\log\frac{1}{p} + (1-p)\log\frac{1}{1-p}$  the binary entropy function.

**Lemma 7** Suppose that  $\beta < 1$ , so that  $R_{d,k} = C(\beta\sigma^2) > C(\sigma^2)$ . For  $\delta > 0$ , denote the corresponding noise level by  $\sigma_0^2(\delta) = C^{-1}((1+\delta)R_{d,k})$  and let  $e(\delta) = \rho_{avg}(C^{-1}((1+\delta)R_{d,k}))$  be the ensemble average decoding error, (20), over the AWGN( $\sigma_0^2$ ) channel. We have that

$$H(\ell|\mathcal{X}_k, \mathbf{Y}) \le \log k - \mathsf{C}(\sigma^2)d + h_b(e(\delta)) + (\delta + e(\delta))\log k$$

and consequently, using (25),

$$I(\mathcal{X}_k; Y) \le h_b(e(\delta)) + (\delta + e(\delta)) \log k. \tag{26}$$

As mentioned above, when  $\beta < 1$ ,  $e(\delta) = o(1)$  for *all* fixed  $\delta > 0$ , by Proposition 4; consequently,  $I(\mathcal{X}_k; \mathbf{Y}) = o(\log k)$ . Combining this with (21), Lemma 5 and (24), assuming sufficiently small  $\varepsilon = O_{\mathbb{R}}(1)$ , we establish (13):

$$\frac{n_{\varepsilon}^*}{\sigma^2 k} \ge C(\varepsilon, \mathsf{R}) \cdot \frac{d}{\sigma^2} \cdot (I(\boldsymbol{\mathcal{X}}_k; \boldsymbol{Y}))^{-1} = \omega \left( \frac{1}{\sigma^2} \cdot \frac{d}{\log k} \right) = \omega \left( \frac{1}{\sigma^2 \mathsf{C}(\beta \sigma^2)} \right) = \omega(1) \,. \tag{27}$$

One can get quantitative bounds by carefully setting  $\delta = o(1)$ , as we do in Appendix, Section B.5:

**Lemma 8** *Suppose that*  $\beta$  < 1. *For small enough fixed*  $\varepsilon \leq \varepsilon_0(R)$ :

1. (Positive rate). If R > 0 then

$$\frac{n_{\varepsilon}^*}{\sigma^2 k} \ge C(\varepsilon, \beta, \mathsf{R}) \sqrt{\frac{\log k}{\log \log k}}. \tag{28}$$

2. (Zero rate). If R = 0 then

$$\frac{n_{\varepsilon}^*}{\sigma^2 k} \ge C(\varepsilon, \beta) \min \left\{ \sqrt{\frac{\log k}{\log \log k}}, \sqrt{\frac{d}{\log k}} \right\}. \tag{29}$$

**Proof** (Of Theorem 2). Directly follows from Lemma 8.

## 4. Proof of Upper Bound

In this section we prove the upper bound of Theorem 1, assuming the rate is smaller than the capacity  $(\beta > 1)$ . The proof is constructive: we propose and analyze an algorithm (which runs in exponential time), whose output  $\hat{\mathcal{X}}_k$  satisfies  $\mathbb{E}\mathcal{L}_{\text{avg}}(\mathcal{X}_k, \hat{\mathcal{X}}_k) \leq \varepsilon$ . It consists of two steps, each using different measurements: Step I is allocated N samples, while Step II uses the remaining  $\bar{N} = n - N$  samples.

Step I consists of a *brute-force search* over an exponential-sized set of candidate centers. Let  $\varepsilon_{\rm I}>0$  be a given precision level, and fix  $\mathcal{T}$  a  $\sqrt{\varepsilon_{\rm I}d/2}$ -net of the sphere  $\sqrt{d}\mathcal{S}^{d-1}$ . For *each* candidate  $\hat{\boldsymbol{X}}\in\mathcal{T}$ , we use the measurements  $\boldsymbol{Y}_1,\ldots,\boldsymbol{Y}_N$  allocated for this step to essentially solve a composite hypothesis testing problem, distinguishing between two alternatives: 1)  $\hat{\boldsymbol{X}}$  is  $\sqrt{\varepsilon d/2}$ -close to *some* center  $\boldsymbol{X}_i$ ; 2)  $\hat{\boldsymbol{X}}$  is  $\sqrt{\varepsilon d}$ -far from *all* the centers. We show that for "typical" center configurations  $\boldsymbol{\mathcal{X}}_k$ , if  $N\gtrsim\sigma^2k\frac{\log(1/\varepsilon_{\rm I})}{\varepsilon_1^2}$  then the test correctly throws away *all* the far points, and retains *most* of the close points. Since the true centers  $\boldsymbol{\mathcal{X}}_k\sim \mathrm{Unif}(\sqrt{d}\mathcal{S}^{d-1})^{\otimes k}$  are (w.h.p.)  $\Omega(\sqrt{d})$ -separated, the remaining points in  $\mathcal{T}$ , that have not been discarded, may be clustered into at most k parts. Step I concludes by returning a list  $\tilde{\boldsymbol{\mathcal{X}}}_{\rm I}$  containing one representation of every cluster.

The dependence of Step I on the precision is sub-optimal: N has to scale like  $\frac{\log(1/\varepsilon_{\rm I})}{\varepsilon_{\rm I}^2}$  instead of  $1/\varepsilon_{\rm I}$ . This sub-optimal rate is mended in Step II. We show that there is a *constant* precision level  $\varepsilon_0$ , that depends on R,  $\beta>1$  (namely, *how much* the rate is smaller than the capacity) so that whenever  $\varepsilon_{\rm I}\leq\varepsilon_0$ , one can construct a mismatched decoder, using  $\tilde{\boldsymbol{\mathcal{X}}}_{\rm I}$ , that *consistently* decodes messages encoded by the true codebook  $\boldsymbol{\mathcal{X}}_k$ . In other words: given a measurement  $\boldsymbol{Y}=\boldsymbol{X}_\ell+\sigma\boldsymbol{Z}$ , one can consistently estimate the unknown label  $\ell$  (up to a global re-labeling). In Step II we observe  $\bar{N}=n-N$  new measurements, and cluster them according to their decoded label. For every cluster

 $i \in [k]$ , we compute the corresponding sample average  $A_i$ , and project it onto the ball  $\mathcal{B}(\mathbf{0}, \sqrt{d})$  to get our final estimate  $\hat{X}_i = \mathcal{P}(A_i)$ . Since each label i witnesses, on average,  $\bar{N}/k$  measurements, the MSE is, to leading order,  $d^{-1}\mathbb{E}||X_i - \hat{X}_i||^2 = \sigma^2 k/\bar{N}$ . Thus, using  $N = C(\mathsf{R}, \beta)\sigma^2 k$  measurements for Step I, and  $\bar{N} = \sigma^2 k/\varepsilon$  Step II, yields a list  $\hat{\mathcal{X}}_k$  with  $\mathbb{E}\mathcal{L}_{avg}(\mathcal{X}_k, \hat{\mathcal{X}}_k) \leq \varepsilon$ .

In the remainder of this section, we provide the full details of the strategy outlined above.

### 4.1. Step I: Brute-Force Search

Let  $\varepsilon_{\rm I}\in(0,1/2)$  a precision parameter. Let  $\mathcal T$  be a fixed  $\sqrt{\varepsilon_{\rm I}d/2}$ -net of  $\sqrt{d}\mathcal S^{d-1}$ , such that  $\forall \pmb X\in\sqrt{d}\mathcal S^{d-1}\exists\hat{\pmb X}\in\mathcal T$  with  $\|\pmb X-\hat{\pmb X}\|^2\leq\varepsilon_{\rm I}d/2$ . By standard estimates, e.g. (Wainwright, 2019, Example 5.8), we can assume that  $|\mathcal T|\leq e^{Cd\log(1/\varepsilon_{\rm I})}$  for some universal C>0. Our goal is to devise a procedure that, given N samples  $\pmb Y_1,\dots,\pmb Y_N\sim {\rm GMM}_{d,k}(\pmb {\mathcal X}_k,\sigma^2)$ , will allow us to discard all candidates  $\hat{\pmb X}\in\mathcal T$  that are  $\sqrt{\varepsilon_{\rm I}d}$ -far from all the centers  $\pmb X_1,\dots,\pmb X_k$ , while keeping enough candidates  $\hat{\pmb X}$  that are  $\sqrt{\varepsilon_{\rm I}d}$ -close to some center; ideally, at least one candidate close to almost every  $\pmb X_i$ . Denote the sets,  $\mathcal H_{\rm Close},\mathcal H_{\rm Far}\subset\mathbb R^d$ 

$$\mathcal{H}_{\text{Close}} = \left\{ \hat{\boldsymbol{X}} : \operatorname{dist}^{2}(\hat{\boldsymbol{X}}, \boldsymbol{\mathcal{X}}_{k}) \leq \frac{1}{2} \varepsilon_{\text{I}} d \right\}, \ \mathcal{H}_{\text{Far}} = \left\{ \hat{\boldsymbol{X}} : \operatorname{dist}^{2}(\hat{\boldsymbol{X}}, \boldsymbol{\mathcal{X}}_{k}) \geq \varepsilon_{\text{I}} d \right\}. \tag{30}$$

We would like a test that, with high probability: 1) rejects all  $\hat{X} \in \mathcal{H}_{Far} \cap \mathcal{T}$ ; 2) accepts most of  $\hat{X} \in \mathcal{H}_{Close} \cap \mathcal{T}$ . Note that since  $\mathcal{T}$  is a  $\sqrt{\varepsilon_I d/2}$ -cover, then for every  $1 \leq i \leq k$ , there is some  $\hat{X} \in \mathcal{H}_{Close} \cap \mathcal{T}$  such that in fact  $\|\hat{X} - X_i\|^2 \leq \varepsilon_I d/2$ .

As a first step, we consider a "local test" Test:  $\mathbb{R}^d \times \mathbb{R}^d \to \{0,1\}$ , that takes a candidate  $\hat{X} \in \mathcal{T}$  and a single sample  $Y \sim \mathrm{GMM}_{d,k}$ , and outputs a decision  $\in \{0,1\}$ . Consider the quantities:

$$q_{\text{Close}}(\boldsymbol{\mathcal{X}}_{k}, \boldsymbol{\mathcal{T}}) = \min_{\hat{\boldsymbol{X}} \in \mathcal{H}_{\text{Close}} \cap \boldsymbol{\mathcal{T}}} \Pr(\mathsf{Test}(\hat{\boldsymbol{X}}, \boldsymbol{Y}) = 1 | \boldsymbol{\mathcal{X}}_{k}),$$

$$q_{\text{Far}}(\boldsymbol{\mathcal{X}}_{k}, \boldsymbol{\mathcal{T}}) = \max_{\hat{\boldsymbol{X}} \in \mathcal{H}_{\text{Far}} \cap \boldsymbol{\mathcal{T}}} \Pr(\mathsf{Test}(\hat{\boldsymbol{X}}, \boldsymbol{Y}) = 1 | \boldsymbol{\mathcal{X}}_{k}).$$
(31)

For a local test Test :  $\mathbb{R}^d \times \mathbb{R}^d \to \{0,1\}$ , a cover  $\mathcal{T}$ , and  $\nu > 0$ , define

$$\mathcal{E}_{\mathsf{Test},\mathcal{T},\nu} = \left\{ \boldsymbol{\mathcal{X}}_k \in \left( \sqrt{d} \mathcal{S}^{d-1} \right)^k : \mathsf{q}_{\mathsf{Close}}(\boldsymbol{\mathcal{X}}_k,\mathcal{T}) \ge \frac{1}{2} k^{-1}, \; \mathsf{q}_{\mathsf{Far}}(\boldsymbol{\mathcal{X}}_k,\mathcal{T}) \le 2k^{-1-\nu} \right\}. \tag{32}$$

We construct a local test with the following properties.

**Lemma 9** Assume that  $\beta > 1$ , and fix a cover  $\mathcal{T}$  of size  $|\mathcal{T}| \leq e^{Cd\log(1/\varepsilon_{\mathrm{I}})}$ . There are positive constants  $\varepsilon_0$ , c, that depend on R,  $\beta$ , and a local test, Test :  $\mathbb{R}^d \times \mathbb{R}^d \to \{0,1\}$  (that depends on  $d,k,\sigma^2$ ) such that for every fixed  $\varepsilon_{\mathrm{I}} \in (0,\varepsilon_0)$ 

$$\lim_{d \to \infty} \Pr\left( \mathcal{X}_k \notin \mathcal{E}_{\mathsf{Test}, \mathcal{T}, c\varepsilon_{\mathrm{I}}^2} \right) = 0.$$
 (33)

We propose a local test Test based on the capacity-achieving decoder used in the proof of Proposition 4. Due to space constraints, the details are deferred to Appendix, Section C.1.

Note that if  $q_{Far}(\mathcal{X}_k, \mathcal{T}) \ll q_{Close}(\mathcal{X}_k, \mathcal{T})$ , as is the case for  $\mathcal{X}_k \in \mathcal{E}_{\mathsf{Test}, \mathcal{T}, \nu}$ , then by observing the statistics of the N local test outputs  $\{\mathsf{Test}(\hat{\boldsymbol{X}}, \boldsymbol{Y}_j)\}_{j=1}^N$ , which are i.i.d. Bernoulli random

variables, one can distinguish between  $\hat{X} \in \mathcal{H}_{\text{Close}}$  and  $\hat{X} \in \mathcal{H}_{\text{Far}}$  with error probability vanishing in N. In particular, consider the candidates  $\hat{X} \in \mathcal{T}$  that pass the following threshold-based test

$$\mathcal{T}_{\text{Close}} = \left\{ \hat{\boldsymbol{X}} \in \mathcal{T} : \sum_{j=1}^{N} \text{Test}(\hat{\boldsymbol{X}}, \boldsymbol{Y}_{j}) \ge \frac{1}{4} k^{-1} N \right\}. \tag{34}$$

**Lemma 10** Fix any  $\mathcal{X}_k \in \mathcal{E}_{\mathsf{Test},\mathcal{T},c\varepsilon_1^2}$  and suppose that

$$N \ge C_1 \sigma^2 k \frac{\log(1/\varepsilon_{\rm I})}{\varepsilon_{\rm I}^2} + C_2 k \log(1/\varphi), \qquad (35)$$

where  $C_1 = C_1(\mathsf{R}, \beta)$ ,  $C_2 > 0$  is a universal constant and  $\varphi \in (0, 1)$ . Then w.p.  $1 - \varphi - o_{\beta, \mathsf{R}}(1)$  over  $Y_1, \ldots, Y_n \sim \mathrm{GMM}_{d,k}(\mathcal{X}_k, \sigma^2)$ , the following event holds:

- 1. (No far candidates).  $\mathcal{T}_{\text{Close}} \cap \mathcal{H}_{\text{Far}} \neq \emptyset$ .
- 2. (Most centers have a cluster). There is  $\mathcal{I} \subseteq [k]$  with  $|\mathcal{I}| \ge (1-\varphi)k$  and  $\max_{i \in \mathcal{I}} \operatorname{dist}^2(\boldsymbol{X}_i, \mathcal{T}_{\operatorname{Close}}) \le \varepsilon_{\operatorname{I}} d$ .

The proof of Lemma 10 appears in the Appendix, Section C.2.

To conclude step I, note that if the minimal distance between centers is  $> 4\sqrt{\varepsilon_I d}$ , then two candidates that are  $\sqrt{\varepsilon_I d}$ -close to different centers are necessarily  $2\sqrt{\varepsilon_I d}$ -far from one another.

Let  $\mathcal{X}_{\mathrm{I}}$  be any  $2\sqrt{\varepsilon_{\mathrm{I}}d}$ -separated subset of  $\mathcal{T}_{\mathrm{Close}}$  of maximal size. We prove in Appendix, Section C.3 that with high probability,  $\mathcal{X}_{k}$  indeed has  $\Omega(\sqrt{d})$  minimal distance, and so:

**Lemma 11** Assume that  $\beta > 1$ ,  $\varepsilon_{\rm I} \leq \varepsilon_0({\sf R},\beta)$  is small enough, and N satisfies (35). W.p.  $1 - \varphi - o_{\beta,{\sf R},\varepsilon_{\rm I}}(1)$  over both  $\mathcal{X}_k \sim {\rm Unif}(\sqrt{d}\mathcal{S}^{d-1})^{\otimes k}$  and  $[\mathbf{Y}_1,\ldots,\mathbf{Y}_n\,|\,\mathcal{X}_k] \sim {\rm GMM}_{d,k}(\mathcal{X}_k,\sigma^2)$ , the following event holds:

- 1.  $\mathcal{X}_{I}$  is a list of size  $(1 \varphi)k \leq m \leq k$ .
- 2. There is  $\mathcal{I} \subseteq [k]$  with  $|\mathcal{I}| = m$  so that for every  $i \in \mathcal{I}$ , there is a unique  $\tilde{\mathbf{X}} \in \tilde{\mathbf{X}}_{\mathrm{I}}$  such that  $\|\mathbf{X}_i \tilde{\mathbf{X}}\|^2 \leq \varepsilon_{\mathrm{I}} d$ .

### 4.2. Step II: Clustering and Averaging

Upon the successful completion of Step I, Lemma 11, we have produced a partial codebook  $\tilde{\mathcal{X}}_{\mathrm{I}}$  of size  $m \geq (1-\varphi)k$ . Moreover, there is a large subset of messages  $\mathcal{I} \subset [k]$ ,  $|\mathcal{I}| = m$  such that for all  $i \in \mathcal{I}$ , the *true*, *unknown* codeword  $X_i$  is  $\sqrt{\varepsilon_1 d}$ -close to a unique codeword  $\tilde{X}_l$  of  $\tilde{\mathcal{X}}_{\mathrm{I}}$ . Provided that  $\varepsilon_{\mathrm{I}} \leq \varepsilon_0(\mathsf{R},\beta)$  is small enough (but constant), it turns out we can construct a "mismatched decoder", using  $\tilde{\mathcal{X}}_{\mathrm{I}}$ , that can consistently decode measurements  $Y = X_\ell + \sigma Z$  in the following sense: 1) If  $\ell \in \mathcal{I}$  then, up to a global relabeling, the decoder returns the correct label  $\ell$ ; 2) If  $\ell \notin \mathcal{I}$ , the decoder consistently returns an error symbol #. Due to space constraints, we defer all the details to Appendix, Section D.1.

In Step II we are given  $\bar{N}=n-N$  new measurements. We use  $\tilde{\boldsymbol{\mathcal{X}}}_{\mathrm{I}}$ , the codebook from Step I, to decode the corresponding labels; measurements for which the decoder returns # are discarded. We end up with m clusters, and for each cluster  $l\in[m]$  we compute the corresponding sample mean  $\boldsymbol{A}_l$ . Finally, we return the list  $\hat{\boldsymbol{\mathcal{X}}}_k=(\hat{\boldsymbol{X}}_1,\ldots,\hat{\boldsymbol{X}}_k)$  such that  $\hat{\boldsymbol{X}}_l=\mathcal{P}(\boldsymbol{A}_l)$  for  $l\in[m],\mathcal{P}(\cdot)$  being the projection onto the ball  $\mathcal{B}(\mathbf{0},\sqrt{d})$ , and  $\hat{\boldsymbol{X}}_l=\mathbf{0}$  for  $m+1\leq l\leq k$ .

The following Lemma bounds the error of the *entire* end-to-end procedure, including both Step I and II. The details are deferred to Appendix, Section D.3.

**Lemma 12** Assume that  $\beta > 1$  and  $\varepsilon \leq \varepsilon_0(R, \beta)$  is small enough. Suppose that

- 1. Step I is run with  $N \ge C\sigma^2 k + Ck \log(1/\varphi)$  measurements,
- 2. Step II is run with  $\bar{N} \geq \frac{k\sigma^2}{\varepsilon} + C\frac{k}{\varepsilon^{1/2}}\log(1/\varphi)$  measurements,

where  $C = C(R, \beta)$  is constant and  $\varphi \in (0, 1)$  is a parameter. Then

$$\lim_{d\to\infty} \mathbb{E}\mathcal{L}_{avg}(\boldsymbol{\mathcal{X}}_k, \hat{\boldsymbol{\mathcal{X}}}_k) \leq \frac{\varepsilon}{1-\varepsilon^{1/4}} + 12\varphi.$$

**Proof** (Of Theorem 1). The claimed lower bound follows from Lemma 6. The upper bound follows by setting, e.g.,  $\varphi = \varepsilon^2$  in Lemma 12, noting that as  $\varepsilon \to 0$ ,  $\frac{\varepsilon}{1-\varepsilon^{1/4}} = \varepsilon + o(\varepsilon)$ .

## Acknowledgments

We are grateful to Uri Erez for helpful discussions. The work of ER and OO is supported in part by the ISF under grant 1641/21. ER is supported in part by an Einstein-Kaye fellowship from the Hebrew University of Jerusalem. TB is supported in part by the ISF grant no. 1924/21, the BSF grant no. 2020159, and the NSF-BSF grant no. 2019752.

## References

Dimitris Achlioptas and Frank McSherry. On spectral learning of mixtures of distributions. In *International Conference on Computational Learning Theory*, pages 458–469. Springer, 2005.

Robert J Adler and Jonathan E Taylor. *Random fields and geometry*. Springer Science & Business Media, 2009.

Animashree Anandkumar, Rong Ge, Daniel Hsu, Sham M Kakade, and Matus Telgarsky. Tensor decompositions for learning latent variable models. *Journal of machine learning research*, 15: 2773–2832, 2014.

Joseph Anderson, Mikhail Belkin, Navin Goyal, Luis Rademacher, and James Voss. The more, the merrier: the blessing of dimensionality for learning large Gaussian mixtures. In *Conference on Learning Theory*, pages 1135–1164. PMLR, 2014.

Sanjeev Arora and Ravi Kannan. Learning mixtures of arbitrary gaussians. In *Proceedings of the thirty-third annual ACM symposium on Theory of computing*, pages 247–257, 2001.

Shiri Artstein-Avidan, Apostolos Giannopoulos, and Vitali D Milman. *Asymptotic geometric analysis, Part I*, volume 202. American Mathematical Soc., 2015.

Hassan Ashtiani, Shai Ben-David, Nicholas JA Harvey, Christopher Liaw, Abbas Mehrabian, and Yaniv Plan. Nearly tight sample complexity bounds for learning mixtures of Gaussians via sample compression schemes. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 3416–3425, 2018.

A. Barg and G.D. Forney. Random codes: minimum distances and error exponents. *IEEE Transactions on Information Theory*, 48(9):2568–2573, 2002.

### ROMANOV BENDORY ORDENTLICH

- Amir Bennatan, A Robert Calderbank, and Shlomo Shamai. Bounds on the MMSE of "bad" LDPC codes at rates above capacity. In 2008 46th Annual Allerton Conference on Communication, Control, and Computing, pages 1065–1072, 2008.
- Aditya Bhaskara, Moses Charikar, and Aravindan Vijayaraghavan. Uniqueness of tensor decompositions with applications to polynomial identifiability. In *Conference on Learning Theory*, pages 742–778. PMLR, 2014.
- Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.
- Ronit Bustin and Shlomo Shamai. MMSE of "bad" codes. *IEEE Transactions on Information Theory*, 59(2):733–743, 2013.
- Siu-On Chan, Ilias Diakonikolas, Rocco A Servedio, and Xiaorui Sun. Efficient density estimation via piecewise polynomial approximation. In *Proceedings of the forty-sixth annual ACM symposium on Theory of computing*, pages 604–613, 2014.
- Thomas M Cover and Joy A Thomas. *Elements of Information Theory*. John Wiley & Sons, 2012.
- Sanjoy Dasgupta. Learning mixtures of Gaussians. In 40th Annual Symposium on Foundations of Computer Science (Cat. No. 99CB37039), pages 634–644. IEEE, 1999.
- Sanjoy Dasgupta and Leonard J Schulman. A probabilistic analysis of EM for mixtures of separated, spherical Gaussians. *Journal of Machine Learning Research*, 8:203–226, 2007.
- Ilias Diakonikolas, Daniel M Kane, and Alistair Stewart. List-decodable robust mean estimation and learning mixtures of spherical Gaussians. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1047–1060, 2018.
- Ilias Diakonikolas, Gautam Kamath, Daniel Kane, Jerry Li, Ankur Moitra, and Alistair Stewart. Robust estimators in high-dimensions without the computational intractability. *SIAM Journal on Computing*, 48(2):742–864, 2019.
- U. Erez and R. Zamir. Achieving 1/2 log (1+SNR) on the AWGN channel with lattice encoding and decoding. *IEEE Transactions on Information Theory*, 50(10):2293–2314, 2004.
- Uri Erez. Lecture notes for principles of coding and detection in communication: Capacity of the AWGN channel. URL https://www.eng.tau.ac.il/~anatolyk/courses/Uri/AWGN\_sphere\_decoder.pdf.
- Jon Feldman, Rocco A Servedio, and Ryan O'Donnell. PAC learning axis-aligned mixtures of Gaussians with no separation assumption. In *International Conference on Computational Learn*ing Theory, pages 20–34. Springer, 2006.
- Robert G Gallager. Information theory and reliable communication, volume 2. Springer, 1968.
- Rong Ge, Qingqing Huang, and Sham M Kakade. Learning mixtures of gaussians in high dimensions. In *Proceedings of the forty-seventh annual ACM symposium on Theory of computing*, pages 761–770, 2015.

- Navin Goyal, Santosh Vempala, and Ying Xiao. Fourier PCA and robust tensor decomposition. In *Proceedings of the forty-sixth annual ACM symposium on Theory of computing*, pages 584–593, 2014.
- Dongning Guo, Shlomo Shamai, and Sergio Verdú. Mutual information and minimum mean-square error in gaussian channels. *IEEE Transactions on Information Theory*, 51(4):1261–1282, 2005.
- Dongning Guo, Shlomo Shamai, and Sergio Verdú. The interplay between information and estimation measures. *Foundations and Trends® in Signal Processing*, 6(4):243–429, 2013.
- Moritz Hardt and Eric Price. Tight bounds for learning a mixture of two gaussians. In *Proceedings* of the forty-seventh annual ACM symposium on Theory of computing, pages 753–760, 2015.
- Samuel B Hopkins and Jerry Li. Mixture models, robustness, and sum of squares proofs. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1021–1034, 2018.
- Daniel Hsu and Sham M Kakade. Learning mixtures of spherical gaussians: moment methods and spectral decompositions. In *Proceedings of the 4th conference on Innovations in Theoretical Computer Science*, pages 11–20, 2013.
- Norman Johnson. Continuous univariate distributions. Wiley, New York, 1994. ISBN 0471584959.
- Grigorii Anatol'evich Kabatiansky and Vladimir Iosifovich Levenshtein. On bounds for packings on a sphere and in space. *Problemy peredachi informatsii*, 14(1):3–25, 1978.
- Adam Tauman Kalai, Ankur Moitra, and Gregory Valiant. Efficiently learning mixtures of two Gaussians. In *Proceedings of the forty-second ACM symposium on Theory of computing*, pages 553–562, 2010.
- Ravindran Kannan, Hadi Salmasian, and Santosh Vempala. The spectral method for general mixture models. *SIAM Journal on Computing*, 38(3):1141–1156, 2008.
- Pravesh K Kothari, Jacob Steinhardt, and David Steurer. Robust moment estimation and improved clustering via sum of squares. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1035–1046, 2018.
- Jeongyeol Kwon and Constantine Caramanis. The EM algorithm gives sample-optimality for learning mixtures of well-separated gaussians. In *Conference on Learning Theory*, pages 2425–2487. PMLR, 2020.
- Jerry Li and Ludwig Schmidt. Robust and proper learning for mixtures of gaussians via systems of polynomial inequalities. In *Conference on Learning Theory*, pages 1302–1382. PMLR, 2017.
- Ankur Moitra and Gregory Valiant. Settling the polynomial learnability of mixtures of gaussians. In 2010 IEEE 51st Annual Symposium on Foundations of Computer Science, pages 93–102. IEEE, 2010.
- Karl Pearson. Contributions to the mathematical theory of evolution. *Philosophical Transactions of the Royal Society of London. A*, 185:71–110, 1894.

### ROMANOV BENDORY ORDENTLICH

- Yury Polyanskiy and Yihong Wu. Lecture notes on information theory. *Lecture Notes for ECE563* (*UIUC*) and, 6(2012-2016):7, 2014.
- Oded Regev and Aravindan Vijayaraghavan. On learning mixtures of well-separated gaussians. In 2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS), pages 85–96. IEEE, 2017.
- Elad Romanov, Tamir Bendory, and Or Ordentlich. Multi-reference alignment in high dimensions: sample complexity and phase transition. *SIAM Journal on Mathematics of Data Science*, 3(2): 494–523, 2021.
- Claude E Shannon. Certain results in coding theory for noisy channels. *Information and control*, 1 (1):6–25, 1957.
- Claude E Shannon. Probability of error for optimal codes in a Gaussian channel. *Bell System Technical Journal*, 38(3):611–656, 1959.
- Nathan Srebro, Gregory Shakhnarovich, and Sam Roweis. An investigation of computational and informational limits in gaussian mixture clustering. In *Proceedings of the 23rd international conference on Machine learning*, pages 865–872, 2006.
- Ananda Theertha Suresh, Alon Orlitsky, Jayadev Acharya, and Ashkan Jafarpour. Near-optimal-sample estimators for spherical gaussian mixtures. *Advances in Neural Information Processing Systems*, 27, 2014.
- Santosh Vempala and Grant Wang. A spectral algorithm for learning mixture models. *Journal of Computer and System Sciences*, 68(4):841–860, 2004.
- Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press, 2019.

## Appendix A. Proof of Proposition 4

As mentioned in the main text, the proof amounts to analyzing a certain sub-optimal decoder for the codebook  $\mathcal{X}_k$ . While the decoders, and their analysis, are not new, we nonetheless provide all the details here as a "warm-up" for things to come.

We consider two different families of decoders, depending on whether one operates in the zero or positive rate regime.

### **A.1. Rate Zero** (R = 0)

**The decoder.** Observe that for a spherical code, the MAP decoder (18) reduces to

$$DecOpt(\mathbf{Y}) = \operatorname*{argmax}_{1 \le i \le k} d^{-1} \langle \mathbf{Y}, \mathbf{X}_i \rangle.$$
 (36)

For the analysis, we consider a sub-optimal decoder, based on thresholding the correlation in (36).

When  $Y = X_i + \sigma Z$ , clearly,  $\mathbb{E}[d^{-1}\langle Y, X_i \rangle] = 1$ , while for  $j \neq i$ ,  $\mathbb{E}[d^{-1}\langle Y, X_j \rangle] = 0$ . Fix thresholds  $0 < \eta_1 \leq \eta_2$ . Consider a decoding rule  $\mathsf{DecCORR}_{\eta_1,\eta_2} : \mathbb{R}^d \to [k] \cup \{\#\}$  so that  $\mathsf{DecCORR}_{\eta_1,\eta_2}(Y) = i$  if and only if both of the following hold:

- 1.  $d^{-1}\langle Y, X_i \rangle \ge 1 \eta_1$ .
- 2. For all  $j \neq i$ ,  $d^{-1}\langle \boldsymbol{Y}, \boldsymbol{X}_i \rangle < 1 \eta_2$ .

Note that since  $\eta_1 \leq \eta_2$ , at most one index  $1 \leq i \leq k$  can satisfy the above. If no such i exists, we set  $\mathsf{DecCORR}_{\eta_1,\eta_2}(Y) = \#$ .

**Analysis.** We proceed to bound the error of the decoder DecCORR $_{\eta_1,\eta_2}$ .

By symmetry of the codebook generating process, the error probability (averaged over the ensemble) does not depend on the particular transmitted message (index). For convenience, throughout this section, we always assume, without loss of generality, that the transmitted message is  $\ell=i$  (and implicitly condition on this event). Thus, the value at the receiver end of the channel is  $\boldsymbol{Y}=\boldsymbol{X}_i+\sigma\boldsymbol{Z}$ .

**Lemma 13** Conditioned on any  $\mathcal{X}_k \in (\sqrt{d}S^{d-1})^k$ ,

$$\Pr\left(d^{-1}\langle \boldsymbol{Y}, \boldsymbol{X}_i \rangle < 1 - \eta_1 \,|\, \boldsymbol{\mathcal{X}}_k\right) \le e^{-\frac{\eta_1^2}{2\sigma^2}d}.$$

**Proof**  $d^{-1}\langle \boldsymbol{X}_i + \sigma \boldsymbol{Z}, \boldsymbol{X}_i \rangle \leq 1 - \eta_1$  is equivalent to  $d^{-1}\langle \sigma \boldsymbol{Z}, \boldsymbol{X}_i \rangle \leq -\eta_1$ . Since  $\boldsymbol{Z} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})$ , we have  $d^{-1}\langle \sigma \boldsymbol{Z}, \boldsymbol{X}_i \rangle \sim \mathcal{N}(0, \sigma^2/d)$ , and the bound follows immediately.

**Lemma 14** For fixed  $\eta_1 > 0$ , define the set

$$\mathbb{X}_i = \left\{ \boldsymbol{\mathcal{X}}_k \in (\sqrt{d}\mathcal{S}^{d-1})^k : \max_{j \in [n] \setminus \{i\}} d^{-1} \langle \boldsymbol{X}_i, \boldsymbol{X}_j \rangle \le \sqrt{\frac{2\log(k-1)}{d} + \frac{\eta_1^2}{\sigma^2}} \right\}.$$

For 
$$\mathcal{X}_k \sim \text{Unif}(\sqrt{d}\mathcal{S}^{d-1})$$
:  $\Pr(\mathcal{X}_k \notin \mathbb{X}_i) \leq e^{-\frac{\eta_1^2}{2\sigma^2}d}$ .

**Proof** By the standard tail bound Lemma 35, for  $t \geq 0$ ,  $\Pr(d^{-1}\langle \boldsymbol{X}_i, \boldsymbol{X}_j \rangle \geq t) \leq e^{-dt^2/2}$ . Taking a union bound over (k-1) choices for  $j \neq i$ ,  $\Pr(\max_{j \neq i} d^{-1}\langle \boldsymbol{X}_i, \boldsymbol{X}_j \rangle \geq t) \leq e^{-dt^2/2 + \log(k-1)}$ . Now set  $t = \sqrt{\frac{2\log(k-1)}{d} + \frac{\eta_1^2}{\sigma^2}}$ .

**Lemma 15** Suppose that  $1 - \eta_2 \ge \sqrt{\frac{2 \log(k-1)}{d} + \frac{\eta_1^2}{\sigma^2}} + \sqrt{\frac{2\sigma^2 \log(k-1)}{d}}$ . Then

$$\Pr\left(\max_{j\neq i} d^{-1}\langle \boldsymbol{Y}, \boldsymbol{X}_{j} \rangle \ge 1 - \eta_{2}\right) \le e^{-\frac{d}{2}\left(1 - \eta_{2} - \sqrt{\frac{2\log(k-1)}{d} + \frac{\eta_{1}^{2}}{\sigma^{2}}} - \sqrt{\frac{2\sigma^{2}\log(k-1)}{d}}\right)^{2}} + e^{-\frac{\eta_{1}^{2}}{2\sigma^{2}}d}. \quad (37)$$

**Proof** Fix  $X_k \in X_i$ , where the set  $X_i$  is from Lemma 14. Writing  $Y = X_i + \sigma Z$ , we note that

$$\Pr\left(\max_{j\neq i} d^{-1}\langle \boldsymbol{Y}, \boldsymbol{X}_j \rangle \geq 1 - \eta_2 \, \big| \, \boldsymbol{\mathcal{X}}_k\right) \leq \Pr\left(\max_{j\neq i} d^{-1}\sigma\langle \boldsymbol{Z}, \boldsymbol{X}_j \rangle \geq 1 - \eta_2 - \sqrt{\frac{2\log(k-1)}{d} + \frac{\eta_1^2}{\sigma^2}} \, \big| \, \boldsymbol{\mathcal{X}}_k\right).$$

Now, each  $d^{-1}\sigma\langle \boldsymbol{Z}, \boldsymbol{X}_j\rangle$  is Gaussian with mean 0 and variance  $\sigma^2/d$ . By a standard bound on the maximum of Gaussian random variables, Lemma 41,  $\mathbb{E}[\max_{j\neq i} d^{-1}\sigma\langle \boldsymbol{Z}, \boldsymbol{X}_j\rangle] \leq \sqrt{\frac{2\sigma^2\log(k-1)}{d}}$ . By the Borell-TIS inequality, Lemma 42, we obtain the first term of (37). The second term is just the bound on  $\Pr(\boldsymbol{\mathcal{X}}_k \notin \mathbb{X}_i)$  from Lemma 14.

**Proof** (Of Proposition 4, case R = 0.)

Combining Lemmas 13 and 15, for every  $\eta_1, \eta_2$  satisfying

$$0 < \eta_1 \le \eta_2 < 1 - \sqrt{\frac{2\log(k-1)}{d} + \frac{\eta_1^2}{\sigma^2}} - \sqrt{\frac{2\sigma^2\log(k-1)}{d}}$$
 (38)

the decoder  $DecCORR_{\eta_1,\eta_2}$  attains average error

$$\Pr\left(i \neq \mathsf{DecCORR}_{\eta_1,\eta_2}(\boldsymbol{X}_i + \sigma \boldsymbol{Z})\right) \leq e^{-\frac{d}{2} \left(1 - \eta_2 - \sqrt{\frac{2\log(k-1)}{d} + \frac{\eta_1^2}{\sigma^2}} - \sqrt{\frac{2\sigma^2\log(k-1)}{d}}\right)^2} + 2e^{-\frac{\eta_1^2}{2\sigma^2}d}. \tag{39}$$

To prove the proposition, it clearly suffices to show that when  $R_{d,k} = C(\beta\sigma^2)$ ,  $\beta > 1$ , then  $\sqrt{\frac{2\log(k-1)}{d} + \frac{\eta_1^2}{\sigma^2}} + \sqrt{\frac{2\sigma^2\log(k-1)}{d}}$  is at most a constant, which is strictly smaller than 1. Indeed, since  $\sigma^2 = \omega(1)$ , the first term is o(1). As for the second term,

$$\sqrt{\frac{2\sigma^2\log(k-1)}{d}} \leq \sqrt{2\sigma^2\mathsf{R}_{d,k}} = \sqrt{2\sigma^2\mathsf{C}(\beta\sigma^2)} \leq \sqrt{1/\beta} < 1\,,$$

where we used  $C(s) = \frac{1}{2} \log(1 + 1/s) \le 1/(2s)$ .

### **A.2. Positive Rate** (R > 0)

**Remark.** The analysis of the previous section (R = 0) unfortunately fails in the positive rate regime, where  $\sigma^2$  is constant. To have any hope of finding  $\eta_1 \le \eta_2$  that satisfy condition (38), it is necessary that (taking  $\eta_1, \eta_2 \to 0$ )

$$\sqrt{\frac{2\log(k-1)}{d}} + \sqrt{\frac{2\sigma^2\log(k-1)}{d}} \le 1.$$

Since  $\log(k-1)/d = \mathsf{R}_{d,k} - O(d^{-1})$ , this constrains the rate as  $\mathsf{R}_{d,k} \leq \frac{1}{\sqrt{2}(1+\sigma)} + O(d^{-1})$ . For small  $\sigma$ , this bound is  $\approx 1/\sqrt{2}$ , while  $\mathsf{C}(\beta\sigma^2) \approx \log(1/\beta\sigma^2)$ . Consequently, for  $\sigma^2 = O(1)$  this condition fails to hold, and the analysis from Section A.1 is not sufficient for proving the existence of capacity-approaching codes. We note that this is a well-known limitation of the *analysis*; specifically, Lemma 15 is too crude. It estimates the maximum over "noise terms",  $\max_{j\neq i} d^{-1}\sigma\langle \boldsymbol{Z}, \boldsymbol{X}_j\rangle$  as if they were all independent. In the zero rate regime, different codewords are essentially orthogonal:  $d^{-1}\langle \boldsymbol{X}_i, \boldsymbol{X}_j\rangle \lesssim \sqrt{\frac{\log k}{d}} = o(1)$ ; consequently, by standard results (e.g. (Adler and Taylor, 2009, 2.2.5)), the maximum is indeed very close to the maximum of i.i.d. Gaussians. When k is exponential, however, this is no longer the case, and the correlations between these noise terms can no longer be neglected once R is sufficiently large. Thus, different techniques are necessary to carry out the analysis (cf. the classical book Gallager (1968)).

The decoder. To overcome the obstruction mentioned above, we consider a different, sub-optimal, decoder, which is similar to Shannon's information density threshold decoder Shannon (1957) for a Gaussian i.i.d. codebook, and to that used in Erez and Zamir (2004), see also Erez and Polyanskiy and Wu (2014). Let  $\alpha = \frac{1}{1+\sigma^2}$  and  $\tau = \sigma^2 \alpha = 1 - \alpha$ . For parameters  $\tau \leq \tau_1 \leq \tau_2$ , consider a decoder DecMMSE $_{\tau_1,\tau_2}: \mathbb{R}^d \to [k] \cup \{\#\}$  so that DecMMSE $_{\tau_1,\tau_2}(Y) = i$  if and only if both of the following hold:

1.  $d^{-1} \|\alpha \mathbf{Y} - \mathbf{X}_i\|^2 \le \tau_1$ . 2. For all  $j \ne i$ ,  $d^{-1} \|\alpha \mathbf{Y} - \mathbf{X}_i\|^2 > \tau_2$ .

If no such  $1 \leq i \leq k$  exists, then  $\mathsf{DecMMSE}_{\tau_1,\tau_2}(\boldsymbol{Y}) = \#.$ 

As was before, in the zero rate case, we analyze the error probability conditioned on the transmitted message being some fixed  $\ell=i$ ; by symmetry, the (ensemble-averaged) error probability does not depend on  $\ell$ . Thus, below,  $\boldsymbol{Y}=\boldsymbol{X}_i+\sigma\boldsymbol{Z}$ .

To justify the name DecMMSE recall that the best linear estimator of  $\boldsymbol{X}_i$  from  $\boldsymbol{Y} = \boldsymbol{X}_i + \sigma \boldsymbol{Z}$ , in the sense of smallest MSE (LMMSE), is  $\alpha \boldsymbol{Y}$ . Note also that  $d^{-1}\mathbb{E}\|\alpha \boldsymbol{Y} - \boldsymbol{X}_i\|^2 = \tau$ , whereas for  $j \neq i$ ,  $d^{-1}\mathbb{E}\|\alpha \boldsymbol{Y} - \boldsymbol{X}_j\|^2 = \alpha^2(1 + \sigma^2) + 1 = \alpha + 1 > \tau$ .

**Analysis.** We proceed to bound the error of the decoder DecMMSE<sub> $\tau_1,\tau_2$ </sub>.

**Lemma 16** For any  $\mathcal{X}_k \in (\sqrt{d}\mathcal{S}^{d-1})^k$ ,

$$\Pr\left(d^{-1}\|\alpha Y - X_i\|^2 > \tau_1 \mid \mathcal{X}_k\right) \le e^{-\frac{1}{2}(1+\sigma^2)(\sqrt{\tau_1/\tau}-1)^2 d}$$

<sup>2.</sup> When  $X_i \sim \mathcal{N}(0,1)^{\otimes d}$  is i.i.d. Gaussian, the LMMSE is actually the MSE-optimal estimator (MMSE). Since we use a spherical prior for  $X_i \sim \mathrm{Unif}(\sqrt{d}\mathcal{S}^{d-1})$ , this is no longer holds exactly, though the discrepancy is negligible when one operates in the regime  $\sigma^2 = \Omega(1)$ .

**Proof** The mapping  $Z \mapsto F(Z) = d^{-1/2} \|\alpha(X_i + \sigma Z) - X_i\|$  is  $d^{-1/2}\alpha\sigma$ -Lipschitz, with expectation

$$\mathbb{E}F(\mathbf{Z}) \leq \sqrt{d^{-1}\mathbb{E}\|\alpha(\mathbf{X}_i + \sigma \mathbf{Z}) - \mathbf{X}_i\|^2} \leq \sqrt{\tau}$$
.

Applying the Gaussian Lipschitz concentration inequality, Lemma 34,

$$\Pr\left( (F(\boldsymbol{Z}))^{2} > \tau_{1} \right) = \Pr\left( F(\boldsymbol{Z}) > \sqrt{\tau_{1}} \right)$$

$$\leq \Pr\left( F(\boldsymbol{Z}) - \mathbb{E}F(\boldsymbol{Z}) > \sqrt{\tau_{1}} - \sqrt{\tau} \, \middle| \, \boldsymbol{\mathcal{X}}_{k} \right)$$

$$\leq e^{-\frac{1}{2} \frac{(\sqrt{\tau_{1}} - \sqrt{\tau})^{2}}{(d^{-1/2}\alpha\sigma)^{2}}}$$

$$= e^{-\frac{\tau}{2\alpha^{2}\sigma^{2}} (\sqrt{\tau_{1}/\tau} - 1)^{2} d}.$$

Now plug  $\tau=\alpha\sigma^2$ ,  $\alpha=1/(1+\sigma^2)$  to get the claimed bound.

**Lemma 17** For  $j \neq i$ ,  $\mathcal{X}_k \sim (\text{Unif}(\sqrt{d}\mathcal{S}^{d-1}))^k$ ,

$$\Pr\left(d^{-1}\|\alpha \mathbf{Y}_i - \mathbf{X}_j\|^2 \le \tau_2\right) \le \left(1 + \frac{1}{\sigma^2}\right)^{1/2} e^{-\left(\mathsf{C}(\sigma^2) - \frac{1}{2}\log\frac{\tau_2}{\tau}\right)d}.$$

**Proof** For a compact convex body  $K \subset \mathbb{R}^d$ , we denote its boundary by  $\partial K$  and surface area by  $\mathrm{Surf}(\partial K)$ . In addition, we denote the Euclidean ball of radius r, centered around  $a \in \mathbb{R}^d$ , by  $\mathcal{B}(a,r)$ .

The event above, whose probability we wish to bound, is equivalent to the event  $X_j \in \mathcal{B}(\alpha Y_i, \sqrt{\tau_2 d})$ . Since  $X_j \sim \mathrm{Unif}(\sqrt{d}\mathcal{S}^{d-1})$ , this probability (conditioned on Y), is given by the surface area ratios  $\mathrm{Surf}(\sqrt{d}\mathcal{S}^{d-1}\cap\mathcal{B}(\alpha Y,\sqrt{\tau_2 d}))/\mathrm{Surf}(\sqrt{d}\mathcal{S}^{d-1})$ . Since  $\partial K\cap L\subseteq \partial(K\cap L)$ , and the surface area of convex sets is monotonic with respect to containment (e.g., (Artstein-Avidan et al., 2015, Theorem B.1.14)),  $\mathrm{Surf}(\sqrt{d}\mathcal{S}^{d-1}\cap\mathcal{B}(\alpha Y,\sqrt{\tau_2 d}))\leq \mathrm{Surf}(\partial\mathcal{B}(\alpha Y,\sqrt{\tau_2 d}))$ . Consequently, the probability is bounded by

$$\frac{\operatorname{Surf}(\partial \mathcal{B}(\alpha Y_i, \sqrt{\tau_2 d}))}{\operatorname{Surf}(\sqrt{d} \mathcal{S}^{d-1})} = \tau_2^{\frac{d-1}{2}} = \tau_2^{-1/2} e^{-\frac{1}{2} \left(\log \frac{1}{\tau} - \log \frac{\tau_2}{\tau}\right) d}.$$

Lastly, use  $\tau_2 \ge \tau$  and  $1/\tau = 1 + 1/\sigma^2$ .

**Proof** (Of Proposition 4, case R > 0.)

Combining Lemmas 16 and 17, along with a union bound over all  $j \neq i$ , the decoding error, averaged over the ensemble  $\mathcal{X}_k \sim \mathrm{Unif}(\sqrt{d}\mathcal{S}^{d-1})^{\otimes k}$ , is bounded as

$$\Pr\left(i \neq \mathsf{DecMMSE}_{\tau_1, \tau_2}(\boldsymbol{Y})\right) \leq e^{-\frac{1}{2}(1+\sigma^2)(\sqrt{\tau_1/\tau}-1)^2d} + (k-1)\left(1 + \frac{1}{\sigma^2}\right)^{1/2}e^{-\left(\mathsf{C}(\sigma^2) - \frac{1}{2}\log\frac{\tau_2}{\tau}\right)d}.$$
(40)

Choose  $\tau_1 = c\tau$ ,  $\tau_2 = c^2\tau$  where c > 1 is a sufficiently small constant. In that case, the first term of (40) clearly decays exponentially in d.

As for the second term, set  $k=e^{d\mathsf{C}(\beta\sigma^2)}$ , so the term is bounded like  $(1+\frac{1}{\sigma^2})^{1/2}e^{-d\left(\mathsf{C}(\sigma^2)-\mathsf{C}(\beta\sigma^2)-\log(c)\right)}$ . Since  $\mathsf{C}(\sigma^2)-\mathsf{C}(\beta\sigma^2)$  is a positive constant, if c>1 is small enough then the term decays exponentially fast in d.

## **Appendix B. Proofs from Section 3**

### B.1. Proof of Lemma 5

We reduce the calculation into a "standard" rate-distortion function (RDF) under MSE distortion.

For technical reasons, it will be more convenient to work with a Gaussian prior on the source signal, rather than the uniform distribution over the sphere. The reason is that the latter distribution is not absolutely continuous with respect to Lebesgue measure (it is supported on a manifold of positive co-dimension, namely,  $\mathcal{S}^{d-1}$ ), so that its differential entropy (in the usual sense) is not well-defined.

Introduce Gaussian random variables,  $\mathcal{G}_k = (G_1, \dots, G_k) \sim \mathcal{N}(0, 1)^{\otimes dk}$ , so that  $X_i = \sqrt{d}G_i / \|G_i\|$ . We have the Markov chain,

$$\mathcal{G}_k \longrightarrow \mathcal{X}_k \longrightarrow \hat{\mathcal{X}}_k$$
, (41)

so by the DPI,

$$I(\mathcal{X}_k; \hat{\mathcal{X}}_k) \ge I(\mathcal{G}_k; \hat{\mathcal{X}}_k).$$
 (42)

The next Lemma shows that if  $\hat{\mathcal{X}}_k$  estimates  $\mathcal{X}_k$  with small distortion, then it also estimates  $\mathcal{G}_k$  with small comparable distortion:

**Lemma 18** Suppose the the Markov chain (41) holds. Then for universal constant  $c_0$ ,

$$\mathbb{E}\mathcal{L}(\boldsymbol{\mathcal{G}}_k, \hat{\boldsymbol{\mathcal{X}}}_k) \leq \left(\sqrt{\mathbb{E}\mathcal{L}(\boldsymbol{\mathcal{X}}_k, \hat{\boldsymbol{\mathcal{X}}}_k)} + c_0 d^{-1/2}\right)^2.$$

The proof of Lemma 18 is straightforward, and deferred to Section B.2.

By assumption,  $\mathbb{E}\mathcal{L}(\mathcal{X}_k, \hat{\mathcal{X}}_k) \leq \varepsilon$  and therefore, by Lemma 18,  $\mathbb{E}\mathcal{L}(\mathcal{G}_k, \hat{\mathcal{X}}_k) \leq (\varepsilon^{1/2} + c_0 d^{-1/2})^2$ . Let  $J = (j_1, \dots, j_k) \in [k]^k$  be indices such that  $j_i \in \operatorname{argmin}_{1 \leq j \leq k} \|\mathcal{G}_i - \hat{X}_j\|^2$ . In other words,  $\mathcal{L}(\mathcal{G}_k, \hat{\mathcal{X}}_k) = (dk)^{-1} \sum_{i=1}^k \|\mathcal{G}_i - \hat{X}_{j_i}\|^2$ . The random variable J is, clearly, deterministic given  $\mathcal{G}_k, \hat{\mathcal{X}}_k$ . By the chain rule for mutual information,

$$I(\mathcal{G}_k; \hat{\mathcal{X}}_k) = I(\mathcal{G}_k; \hat{\mathcal{X}}_k, J) - I(\mathcal{G}_k; J | \hat{\mathcal{X}}_k).$$
(43)

Since J is a discrete random variable,

$$I(\boldsymbol{\mathcal{G}}_k; J|\hat{\boldsymbol{\mathcal{X}}}_k) := H(J|\hat{\boldsymbol{\mathcal{X}}}_k) - H(J|\boldsymbol{\mathcal{G}}_k, \hat{\boldsymbol{\mathcal{X}}}_k) \le H(J) \le \log(k^k) = k \log k,$$
(44)

where we used the standard facts that the entropy of a discrete variable is non-negative, and that conditioning decreases entropy.

Set  $D_i = X_{j_i}$  and  $\mathcal{D} = (D_1, \dots, D_k) \in \mathbb{R}^{k \times d}$ , so that, by definition,

$$(dk)^{-1}\mathbb{E}\|\mathcal{G}_k - \mathcal{D}\|_F^2 = \mathbb{E}\mathcal{L}_{avg}(\mathcal{G}_k, \mathcal{D}) \leq (\varepsilon^{1/2} + c_0 d^{-1/2})^2$$
.

Since  $\mathcal{D}$  is a function of  $(\hat{\mathcal{X}}_k, J)$ , the DPI implies  $I(\mathcal{G}_k; \mathcal{D}) \leq I(\mathcal{G}_k; \hat{\mathcal{X}}_k, J)$ . Thus,

$$I(\boldsymbol{\mathcal{G}}_{k}; \hat{\boldsymbol{\mathcal{X}}}_{k}, J) \geq I(\boldsymbol{\mathcal{G}}_{k}; \boldsymbol{\mathcal{D}})$$

$$\geq \min_{P_{\tilde{\boldsymbol{\mathcal{D}}}|\boldsymbol{\mathcal{G}}_{k}}: (dk)^{-1} \|\boldsymbol{\mathcal{G}}_{k} - \tilde{\boldsymbol{\mathcal{D}}}\|^{2} \leq (\varepsilon^{1/2} + c_{0}d^{-1/2})^{2}} I(\boldsymbol{\mathcal{G}}_{k}; \tilde{\boldsymbol{\mathcal{D}}})$$

$$= \frac{dk}{2} \log \left( \frac{1}{(\varepsilon^{1/2} + c_{0}d^{-1/2})^{2}} \right), \tag{45}$$

where (45) is the solution to the classical Gaussian source rate-distortion problem (Polyanskiy and Wu, 2014, Chapter 27). The proof of Lemma 5 concludes by combining (42)-(45).

We remark that for sufficiently small  $\varepsilon$  the lower bound

$$\min_{P_{\hat{\boldsymbol{\mathcal{X}}}_{k}|\boldsymbol{\mathcal{G}}_{k}}:\mathbb{E}\mathcal{L}_{\mathbf{avg}}(\boldsymbol{\mathcal{G}}_{k},\boldsymbol{\mathcal{D}})\leq\varepsilon}I(\boldsymbol{\mathcal{G}}_{k};\hat{\boldsymbol{\mathcal{X}}}_{k})\geq\frac{dk}{2}\log\left(\frac{1}{\varepsilon}\right)-k\log k,\tag{46}$$

which we derived within the proof above, is in fact tight (up to the difference between  $k \log k$  and  $\log |\mathcal{S}_k|$ , where  $\mathcal{S}_k$  is the symmetric group of permutations on [k]). To see this, we consider the Markov chain  $\mathcal{G}_k \to \hat{\mathcal{G}}_k \to \hat{\mathcal{X}}_k$ , where the channel from  $\mathcal{G}_k$  to  $\hat{\mathcal{G}}_k$  is the test channel attaining the Gaussian RDF (see e.g., (Cover and Thomas, 2012, Theorem 10.3.2)), and the channel from  $\hat{\mathcal{G}}_k \to \hat{\mathcal{X}}_k$  is defined by applying a uniform random permutation J on  $\hat{\mathcal{G}}_k$ , resulting in  $\hat{\mathcal{X}}_k$ . Note that

$$I(\mathcal{G}_k; \hat{\mathcal{X}}_k, J) = I(\mathcal{G}_k; \hat{\mathcal{G}}_k) = \frac{dk}{2} \log \left(\frac{1}{\varepsilon}\right)$$
(47)

and that

$$I(\boldsymbol{\mathcal{G}}_k; J|\hat{\boldsymbol{\mathcal{X}}}_k) = H(J|\hat{\boldsymbol{\mathcal{X}}}_k) - H(J|\boldsymbol{\mathcal{G}}_k, \hat{\boldsymbol{\mathcal{X}}}_k) = H(J) - H(J|\boldsymbol{\mathcal{G}}_k, \hat{\boldsymbol{\mathcal{X}}}_k) \approx H(J), \tag{48}$$

where the last approximation is due to the fact that for small  $\varepsilon$  we can recover J from  $\mathcal{G}_k$  and  $\hat{\mathcal{X}}_k$ . Thus, the approximate tightness of (46) follows from (43).

The subtractive  $k \log k$  term we lose here is the reason that the lower bound in Theorem 1 is  $e^{-2R}$  instead of 1. While we believe that 1 is the correct lower bound, this loss seems to be inherent to the mutual information bounding program we follow here.

### **B.2. Proof of Lemma 18**

It is a well-known fact (Johnson, 1994, Eq. 18.15) that  $\mathbb{E}\|G_i\| = \sqrt{d} + O(d^{-1/2})$ . Consequently,

$$\mathbb{E}\|\boldsymbol{G}_i - \boldsymbol{X}_i\|^2 = 2d - 2\sqrt{d}\mathbb{E}\|\boldsymbol{G}_i\| = O(1).$$

Let  $J = (j_1, \dots, j_k)$  be  $j_i = \operatorname{argmin}_{1 < j < k} \| \boldsymbol{X}_i - \hat{\boldsymbol{X}}_i \|$ . By definition of  $\mathcal{L}_{avg}(\cdot, \cdot)$ , (5),

$$\mathcal{L}_{\mathsf{avg}}(\boldsymbol{\mathcal{G}}_k, \hat{\boldsymbol{\mathcal{X}}}_k) \leq rac{1}{dk} \sum_{i=1}^k \mathbb{E} \| \boldsymbol{G}_i - \hat{\boldsymbol{X}}_{j_i} \|^2,$$

while  $\mathcal{L}_{\text{avg}}(\boldsymbol{\mathcal{X}}_k, \hat{\boldsymbol{\mathcal{X}}}_k) = \frac{1}{dk} \sum_{i=1}^k \mathbb{E} \|\boldsymbol{X}_i - \hat{\boldsymbol{X}}_{j_i}\|^2$ . Moreover, observe that

$$oldsymbol{\mathcal{D}} = (oldsymbol{D}_1, \dots, oldsymbol{D}_k) \mapsto \left(rac{1}{dk} \sum_{i=1}^k \mathbb{E} \|oldsymbol{D}_i\|^2
ight)^{1/2}$$

defines a semi-norm on  $d \times k$  matrices (with square-integrable entries). Thus, by the triangle inequality,

$$\begin{split} \left(\mathbb{E}\mathcal{L}_{\text{avg}}(\boldsymbol{\mathcal{G}}_{k}, \hat{\boldsymbol{\mathcal{X}}}_{k})\right)^{1/2} &\leq \left(\frac{1}{dk} \sum_{i=1}^{k} \mathbb{E}\|\boldsymbol{G}_{i} - \hat{\boldsymbol{X}}_{j_{i}}\|^{2}\right)^{1/2} \\ &= \left(\frac{1}{dk} \sum_{i=1}^{k} \mathbb{E}\|(\boldsymbol{G}_{i} - \boldsymbol{X}_{i}) + (\boldsymbol{X}_{i} - \hat{\boldsymbol{X}}_{j_{i}})\|^{2}\right)^{1/2} \\ &\leq \underbrace{\left(\frac{1}{dk} \sum_{i=1}^{k} \mathbb{E}\|\boldsymbol{G}_{i} - \boldsymbol{X}_{i}\|^{2}\right)^{1/2}}_{\left(\frac{1}{dk} \sum_{i=1}^{k} O(1)\right)^{1/2} = O(d^{-1/2})} + \underbrace{\left(\frac{1}{dk} \sum_{i=1}^{k} \mathbb{E}\|\boldsymbol{X}_{i} - \hat{\boldsymbol{X}}_{j_{i}}\|^{2}\right)^{1/2}}_{\left(\mathcal{L}_{\text{avg}}(\boldsymbol{\mathcal{X}}_{k}, \hat{\boldsymbol{\mathcal{X}}}_{k})\right)^{1/2}}. \end{split}$$

### B.3. Proof of Lemma 6

Write  $I(\boldsymbol{\mathcal{X}}_k;\boldsymbol{Y}_1,\ldots,\boldsymbol{Y}_n,\boldsymbol{\ell})=I(\boldsymbol{\mathcal{X}}_k;\boldsymbol{\ell})+I(\boldsymbol{\mathcal{X}}_k;\boldsymbol{Y}_1,\ldots,\boldsymbol{Y}_n|\boldsymbol{\ell})$ , with  $I(\boldsymbol{\mathcal{X}}_k;\boldsymbol{\ell})=0$ . For  $1\leq i\leq k$ , let  $n_i=n_i(\boldsymbol{\ell})=|\boldsymbol{\ell}^{-1}(i)|$  be the number of measurements labeled i. The proof amounts to the following observation: the desired MI  $I(\boldsymbol{\mathcal{X}}_k;\boldsymbol{Y}_1,\ldots,\boldsymbol{Y}_n|\boldsymbol{\ell})$  is simply the cumulative MI across k parallel Gaussians channel, with independent inputs  $\boldsymbol{X}_1,\ldots,\boldsymbol{X}_k$ , such that one observes  $n_i$  outputs (samples) of each channel i. We now quantify this statement.

Let  $I(\sigma^2, m) = I(X; X + \sigma Z_1, ..., X + \sigma Z_m)$  be the input-output MI between X and m outputs through an AWGN( $\sigma^2$ ) channel. Since the sample mean is a sufficient statistic for the true mean under a Gaussian measurement model, we have

$$I(\sigma^2, m) = I\left(\boldsymbol{X}, \frac{1}{m}\left((\boldsymbol{X} + \sigma \boldsymbol{Z}_1, \dots, \boldsymbol{X} + \sigma \boldsymbol{Z}_m)\right)\right) = I\left(\frac{\sigma^2}{m}, 1\right) \le d\mathsf{C}(\sigma^2/m),$$

where  $C(\cdot)$  denotes the AWGN channel capacity (9). Thus,

$$I(\boldsymbol{\mathcal{X}}_k; \boldsymbol{Y}_1, \dots, \boldsymbol{Y}_n | \boldsymbol{\ell}) = \mathbb{E}\left[\sum_{i=1}^k I(\sigma^2, n_i(\boldsymbol{\ell}))\right] \leq \mathbb{E}\left[\sum_{i=1}^k d\mathsf{C}(\sigma^2/n_i(\boldsymbol{\ell}))\right].$$

One may readily verify that the function  $m \mapsto C(\sigma^2/m)$  is concave. By Jensen's inequality,

$$\sum_{i=1}^k d\mathsf{C}(\sigma^2/n(\boldsymbol{\ell})) = k \cdot \frac{1}{k} \sum_{i=1}^k d\mathsf{C}(\sigma^2/n(\boldsymbol{\ell})) \leq k d\mathsf{C}\left(\frac{\sigma^2}{\frac{1}{k} \sum_{i=1}^k n_i(\boldsymbol{\ell})}\right) = k d\mathsf{C}\left(\sigma^2 k/n\right) \,,$$

and the claimed result follows.

We remark that to prove the bound, we *did not* actually need to use the fact that the labels all have the same probability; the calculation above shows that a balanced label distribution in fact maximizes the MI between  $\mathcal{X}_k$  and the observations  $Y_1, \ldots, Y_n$  (though this will not be used later).

### **B.4. Proof of Lemma 7**

The proof relies on the celebrated I-MMSE relation of Guo, Shamai and Verdu (see Guo et al. (2005), Guo et al. (2013), and also the works Bennatan et al. (2008) and Bustin and Shamai (2013) that apply the I-MMSE framework for studying the MSE of estimating the transmitted codeword from the output of the AWGN channel).

Let  $\ell \sim \operatorname{Unif}([k])$  and  $\boldsymbol{Y}(s) = \boldsymbol{X}_{\ell} + \sqrt{s}\boldsymbol{Z}$ . Denote  $I(s) = I(\ell; \boldsymbol{Y}(s) | \boldsymbol{\mathcal{X}}_k) = I(\boldsymbol{X}_{\ell}; \boldsymbol{Y}(s) | \boldsymbol{\mathcal{X}}_k)$ , where equality holds since, with probability one,  $\boldsymbol{X}_1, \dots, \boldsymbol{X}_k$  are all distinct. Recall that our goal is to bound  $H(\ell|\boldsymbol{Y}(\sigma^2), \boldsymbol{\mathcal{X}}_k) = H(\ell|\boldsymbol{\mathcal{X}}_k) - I(\sigma^2)$ .

Clearly, for any  $\sigma_0^2 < \sigma^2$ ,

$$H(\ell|\boldsymbol{Y}(\sigma^2),\boldsymbol{\mathcal{X}}_k) - H(\ell|\boldsymbol{Y}(\sigma_0^2),\boldsymbol{\mathcal{X}}_k) = I(\sigma_0^2) - I(\sigma^2) = -\int_{\sigma_0^2}^{\sigma^2} \frac{d}{ds} I(s) ds.$$

Using the I-MMSE relation, Lemma 44, applied pointwise conditioned on  $\mathcal{X}_k$ ,

$$\frac{d}{ds}I(s) = -\frac{1}{2s^2}\mathbb{E}\left[\left\|\boldsymbol{X}_{\ell} - \mathbb{E}(\boldsymbol{X}_{\ell}|\boldsymbol{Y}(s),\boldsymbol{\mathcal{X}}_k)\right\|^2\right].$$

Since  $\mathbb{E}(\boldsymbol{X}_{\ell}|\boldsymbol{Y}(s),\boldsymbol{\mathcal{X}}_{k})$  is the minimum MSE estimator of  $\boldsymbol{X}_{\ell}$  from  $(\boldsymbol{Y}(s),\boldsymbol{\mathcal{X}}_{k})$ , it holds that for any  $(\boldsymbol{Y}(s),\boldsymbol{\mathcal{X}}_{k})$ -measurable random variable  $\hat{\boldsymbol{X}}=\hat{\boldsymbol{X}}(\boldsymbol{Y}(s),\boldsymbol{\mathcal{X}}_{k})$ , we have  $-\frac{d}{ds}I(s)\leq \frac{1}{2s^{2}}\mathbb{E}\|\boldsymbol{X}_{\ell}-\hat{\boldsymbol{X}}\|^{2}$ .

Choose the optimal linear estimator (LMMSE) of  $X_\ell$  from Y(s), namely  $\hat{X} = \alpha(s)Y(s)$ ,  $\alpha(s) = \frac{1}{1+s}$ , so that  $\mathbb{E}\|X_\ell - \hat{X}\|^2 = \frac{s}{1+s}d$ . One would think, at first sight, that this upper bound should be very loose: after all, the LMMSE is optimal for a Gaussian signal, whereas, conditioned on  $\mathcal{X}_k$ , the distribution of  $X_\ell$  is very much non-Gaussian; it is not even continuous! Recall, however, that we are interested in applying Lemma 7 when the rate  $R_{d,k}$  is *above* the capacity  $C(\sigma^2)$ ; the key intuition is that when this is the case, the joint statistics of  $(X_\ell, Y = X_\ell + \sigma Z)$  with  $\ell \sim \mathrm{Unif}([k])$ , are in some sense "indistinguishable" from those of a joint Gaussian distribution  $(W, Y = W + \sigma Z)$ ,  $W \sim \mathcal{N}(0, I)$ , corresponding to the capacity-achieving distribution of the Gaussian channel.

Continuing the calculation,

$$-\int_{\sigma_0^2}^{\sigma^2}\frac{d}{ds}I(s)ds \leq \int_{\sigma_0^2}^{\sigma^2}\frac{1}{2s^2}\cdot\frac{s}{1+s}ds\cdot d = \int_{\sigma_0^2}^{\sigma^2}\left(-\mathsf{C}'(s)\right)ds\cdot d = \left(\mathsf{C}(\sigma_0^2)-\mathsf{C}(\sigma^2)\right)d\,,$$

where  $C(s) = \frac{1}{2} \log(1 + 1/s)$  is from (9) and C'(s) is its derivative. Combining,

$$H(\ell|\boldsymbol{Y}(\sigma^2), \boldsymbol{\mathcal{X}}_k) \le H(\ell|\boldsymbol{Y}(\sigma_0^2), \boldsymbol{\mathcal{X}}_k) - \mathsf{C}(\sigma^2)d + \mathsf{C}(\sigma_0^2)d. \tag{49}$$

Now, set  $\sigma_0^2 = \mathsf{C}^{-1}((1+\delta)\mathsf{R}_{d,k})$ . To apply (49), we need to verify that  $\sigma_0^2 < \sigma^2$ . Applying the decreasing function  $\mathsf{C}(\cdot)$ , the condition is equivalent to  $\mathsf{C}(\sigma^2) < \mathsf{C}(\sigma_0^2) = (1+\delta)\mathsf{R}_{d,k}$ , which certainly hold since we assume  $\mathsf{R}_{d,k} > \mathsf{C}(\sigma^2)$ .

By definition,  $C(\sigma_0^2)d = (1 + \delta) \log k$ .

Define by  $e(\delta)$  the error (averaged over the ensemble  $\mathcal{X}_k$ ) for decoding  $\ell$  under AWGN( $\sigma_0^2$ ), using codebook  $\mathcal{X}_k$ . In other words, it is the error of the MAP estimator for  $\ell$  given  $(\mathbf{Y}(\sigma_0^2), \mathcal{X}_k)$ . By Fano's inequality, Lemma 43,  $H(\ell|\mathbf{Y}(\sigma_0^2), \mathcal{X}_k) \leq h_b(e(\delta)) + e(\delta) \log k$ . Combined with (49), we obtain the bound claimed in Lemma 7.

### **B.5. Proof of Lemma 8**

Before getting to the computation, we emphasize that Lemma 7 may be invoked with *any* other upper bound on the ensemble average error  $\rho_{\rm avg}(\cdot)$ , that could possibly be obtained through other means, e.g., by analyzing a different decoder than the one from Section A. There is much literature devoted to computing optimal error rates for both the Gaussian i.i.d. and the spherical code ensembles, primary in the regime of *positive rate*. In particular, for rates between the so-called *critical rate* and capacity the exact exponential decay rate is known:  $\rho_{\rm avg}(\sigma^2) = \exp(-E_{SP}^*({\sf R},\sigma^2)d + o(d))$ , where  $E_{SP}^*({\sf R},\sigma^2)$  is the *sphere-packing error exponent*. See, for example, Shannon (1959); Gallager (1968) for the exact expression. In the analysis that follows, we will need bounds on the error probability in the regime  $C(\sigma^2) - R = o(1)$ . In fact, for the zero rate regime, the capacity itself is o(1), and sometimes it decays even as  $o(d^{-1/2})$ . In those cases, the sphere packing error exponent is of limited use.

Instead, we use the upper bounds on  $\rho_{avg}(\cdot)$  derived in Section A.

As before, the analysis is divided between the positive (R > 0) and zero (R = 0) rate regimes.

### **B.5.1. POSITIVE RATE**

Let us work under the slightly more general regime, where  $R_{d,k}$  is either positive or decays slow enough with d, specifically,  $R_{d,k} = \frac{\log k}{d} = \omega(d^{-1/2})$  as  $d \to \infty$ .

We apply the bound (40) with noise variance

$$\sigma_0^2 = \mathsf{C}^{-1}((1+\delta)\mathsf{R}_{d,k})$$
.

The second term of (40) is bounded By

$$(1+1/\sigma_0^2)^{1/2}e^{-d\left((1+\delta)\mathsf{R}_{d,k}-\mathsf{R}_{d,k}-\frac{1}{2}\log(\tau_2/\tau)\right)}=O(1)\cdot e^{-d(\delta\mathsf{R}_{d,k}-\frac{1}{2}\log(\tau_2/\tau))}.$$

Set  $\tau_1 = \tau_2 = (1 + \frac{1}{2} R_{d,k} \delta) \tau$ , so that  $\log(\tau_2/\tau) \le \frac{1}{2} R_{d,k} \delta$ . Thus,

$$O(1) \cdot e^{-d(\delta R_{d,k} - \frac{1}{2}\log(\tau_2/\tau))} = O(1) \cdot e^{-d(\delta R_{d,k} - \frac{1}{4}R_{d,k}\delta)} \le e^{-CR_{d,k}\delta d}.$$

for some C > 0. On the other hand, the first term of (40) is

$$e^{-\frac{1}{2}(1+\sigma_0^2)(\sqrt{\tau_1/\tau}-1)^2d} \lesssim e^{-CR_{d,k}^2\delta^2d} = e^{-C\delta^2\frac{(\log k)^2}{d}}.$$

Note that since  $\delta = o(1)$ , this term is the most significant.

Denote  $A = \frac{(\log k)^2}{d}$ ; recall that for  $R_{d,k} = \omega(d^{-1/2})$ ,  $A = \omega(1)$ .

In light of the estimates above, we need to choose  $\delta = o(1)$  so to minimize  $\delta + e^{-C\delta^2 A}$ . Take

$$\delta = C_1 \sqrt{\frac{\log A}{A}} \,,$$

for large enough constant  $C_1$ , which yields

$$\delta + e(\delta) \lesssim \delta + e^{-C\delta^2 A} \lesssim \sqrt{\frac{\log A}{A}}$$
.

Plugging this into (26),

$$I(\boldsymbol{\mathcal{X}}_k; \boldsymbol{Y}) \lesssim \log k \sqrt{\frac{\log A}{A}}$$
.

Using (27),

$$\begin{split} \frac{n_{\varepsilon}^*}{\sigma^2 k} &\geq C_1(\varepsilon, \mathsf{R}) \cdot \frac{d}{\sigma^2} \cdot (I(\boldsymbol{\mathcal{X}}_k; \boldsymbol{Y}))^{-1} \\ & \gtrsim \frac{d}{\sigma^2 \log k} \sqrt{\frac{A}{\log A}} \end{split}$$

where  $A=\frac{(\log k)^2}{d}$ . Let us understand the asymptotic of this bound as  $d\to\infty$  and  $\mathsf{R}_{d,k}=\mathsf{C}(\beta\sigma^2)=\frac{\log k}{d}\gg d^{-1/2}$ . In that case,  $\log A\approx\log\log k$ , and so, the above reads

$$\frac{n_{\varepsilon}^*}{\sigma^2 k} \gtrsim \frac{d}{\sigma^2 \log k} \sqrt{\frac{A}{\log A}} \gtrsim \frac{d}{\sigma^2 \log k} \sqrt{\frac{(\log k)^2}{d} \log \log k}}.$$

Using  $1/(\beta\sigma^2) \ge \mathsf{C}(\beta\sigma^2) = (\log k)/d$  (since  $\mathsf{C}(s) \le 1/(2s)$ ) finally yields

$$\frac{n_{\varepsilon}^*}{\sigma^2 k} \gtrsim \sqrt{\frac{\log k}{d}} \sqrt{\frac{\log k}{\log \log k}} \,. \tag{50}$$

Finally, note that in the positive rate regime,  $\frac{\log k}{d} = \Omega(1)$ .

## B.5.2. RATE ZERO (R = 0)

Assume that  $\lim_{d\to\infty} \mathsf{R}_{d,k}=0$  (including, possibly,  $\mathsf{R}_{d,k}\gg d^{-1/2}$ ). We would like to use the bound (39) with some  $\eta=\eta_1=\eta_2=o(1)$  and  $\sigma_0^2=\sigma_0^2(\delta)=0$  $C^{-1}((1+\delta)R_{d,k})$ , for  $\delta = o(1)$ .

We start with the condition (38), namely,

$$0 \le F := 1 - \eta - \sqrt{\frac{2\log k}{d} + \frac{\eta^2}{\sigma_0^2(\delta)}} - \sqrt{\frac{2\sigma_0^2(\delta)\log k}{d}}$$
$$= 1 - \eta - \sqrt{2\mathsf{R}_{d,k}}\sqrt{1 + \frac{1}{2\sigma_0^2(\delta)\mathsf{R}_{d,k}}\eta^2} - \sqrt{2\sigma_0^2(\delta)\mathsf{R}_{d,k}}.$$

(we replace k-1 with k, which yields a stronger condition.)

Use  $C(s) = \frac{1}{2} \log(1 + 1/s) \le 1/(2s)$ , therefore  $C^{-1}(y) \le 1/(2y)$ , and so  $\sigma_0^2(\delta) \le \frac{1}{2(1+\delta)R_{JL}}$ :

$$\sqrt{2\sigma_0^2(\delta)\mathsf{R}_{d,k}} \leq \sqrt{\frac{1}{1+\delta}} = 1 - \frac{\delta}{2} + O(\delta^2) \,.$$

Moreover,

$$\sqrt{1 + \frac{1}{2\sigma_0^2(\delta)\mathsf{R}_{d,k}}\eta^2} \le 1 + \frac{1}{4\sigma_0^2(\delta)\mathsf{R}_{d,k}}\eta^2$$

 $(\sqrt{1+x} \le 1 + \frac{1}{2}x \text{ for all } x \ge 0).$  Thus,

$$F \ge \delta/2 - O(\delta^2) - \eta - \sqrt{2R_{d,k}} - 2\sqrt{R_{d,k}} \cdot \frac{1}{4\sigma_0^2(\delta)R_{d,k}} \eta^2.$$
 (51)

Note that  $\sigma_0^2(\delta)\mathsf{R}_{d,k} = \Theta(1)$  for any  $\delta = o(1)$ ; to see this, recall that  $(1+\delta)\mathsf{R}_{d,k} = \mathsf{C}(\sigma_0^2)$  (by definition), with  $\mathsf{C}(\sigma_0^2) = 1/(2\sigma_0^2) + O(1/\sigma_0^4)$  with  $\sigma_0^2 \to \infty$ . Consequently, the last term of (51) above is necessarily of lower order than either  $\sqrt{\mathsf{R}_{d,k}}$  or  $\eta$ .

Introduce a constant parameter  $\nu \in (0, 1/2)$ , and set  $\eta = (1/2 - \nu)\delta$ . Observe that whenever

$$\delta \ge \frac{2\sqrt{2}}{\nu} \sqrt{\mathsf{R}_{d,k}} = \frac{2\sqrt{2}}{\nu} \sqrt{\frac{\log k}{d}}, \quad \delta = o(1),$$
 (52)

plugging into (51), we have  $F \ge \frac{\nu}{2}\delta(1 - o(1)) > 0$ .

Let us estimate the terms in (39). The first term is  $e^{-\frac{d}{2}F^2} \leq e^{-C_1 d\nu^2 \delta^2}$ . The second term is  $2e^{-\frac{\eta^2}{2}\cdot\frac{d}{\sigma_0^2}}$ . Using

$$\frac{\eta^2 d}{2\sigma_0^2} \ge \eta^2 d(1+\delta) \mathsf{R}_{d,k} = \eta^2 (1+\delta) \log k \ge \eta^2 \log k \,,$$

(we used  $\sigma_0^2 \leq \frac{1}{2(1+\delta)\mathsf{R}_{d,k}}$ ), we deduce that the second term is  $\leq 2e^{-(1/2-\nu)^2\delta^2\log k}$ . Since in the zero rate regime,  $d\gg\log k$ , we see that the first term is always negligible compared to the second, regardless of how fast  $\delta$  decays. Thus, we would like choose  $\delta=o(1)$  so to minimize (the asymptotic decay rate of)

$$\delta + e(\delta) \lesssim \underbrace{\delta}_{e_1(\delta)} + \underbrace{e^{-(1/2 - \nu)^2 \delta^2 \log k}}_{e_2(\delta)} . \tag{53}$$

Note that  $e_1(\delta)$  is increasing in  $\delta$ , while  $e_2(\delta)$  is decreasing. Denote

$$\delta_1 = \frac{\sqrt{2}}{(1/2 - \nu)} \sqrt{\frac{\log \log k}{\log k}}, \quad \delta_2 = \frac{2\sqrt{2}}{\nu} \sqrt{\frac{\log k}{d}}, \tag{54}$$

so that  $\delta_2$  is the smallest number  $\delta$  that satisfies (52).

One may readily verify that  $\delta = \delta_1$  optimally balances between  $e_1(\delta), e_2(\delta)$ , in the sense of asymptotic growth:

$$e_1(\delta_1) \approx e_2(\delta_1) \approx \sqrt{\frac{\log \log k}{\log k}}, \implies e_1(\delta_1) + e_2(\delta_1) \lesssim \sqrt{\frac{\log \log k}{\log k}}$$

Recall, however, that not all assignments  $\delta$  are applicable; we must satisfy the constraint (52),  $\delta \geq \delta_2$ . If  $\delta_2 \leq \delta_1$  then there is no problem; on the other hand, if  $\delta_2 > \delta_1$ , assigning  $\delta = \delta_2$ ,

$$e_1(\delta_2) + e_2(\delta_2) \stackrel{(i)}{\leq} e_1(\delta_2) + e_2(\delta_1) \stackrel{(ii)}{\lesssim} e_1(\delta_2) + e_1(\delta_1) \stackrel{(iii)}{\leq} 2e_1(\delta_2) \lesssim \sqrt{\frac{\log k}{d}},$$

where we used that: (i)  $e_2(\cdot)$  is decreasing; (ii)  $e_1(\delta_1) \approx e_2(\delta_1)$ ; (iii)  $e_1(\cdot)$  is increasing.

Concluding the calculation, using (26), we have

$$I(\mathcal{X}_k; \mathbf{Y}) \lesssim \log k \cdot \max \left\{ \sqrt{\frac{\log \log k}{\log k}}, \sqrt{\frac{\log k}{d}} \right\}.$$
 (55)

Finally, to deduce the lower bound on the sample complexity, use (27):

$$\frac{n_{\varepsilon}^{*}}{\sigma^{2}k} \geq C_{1}(\varepsilon) \cdot \frac{d}{\sigma^{2}} \cdot (I(\boldsymbol{\mathcal{X}}_{k}; \boldsymbol{Y}))^{-1} 
\geq C_{1}(\varepsilon) \cdot 2\beta \log k \cdot (I(\boldsymbol{\mathcal{X}}_{k}; \boldsymbol{Y}))^{-1} 
\geq C_{2}(\varepsilon)\beta \min \left\{ \sqrt{\frac{\log k}{\log \log k}}, \sqrt{\frac{d}{\log k}} \right\}.$$
(56)

28

## Appendix C. Proofs for Section 4.1

### C.1. Proof of Lemma 9

As in Section A, we give different constructions between the zero rate (R = 0) and positive rate (R > 0) regimes. The construction for the local test is guided by the form of the capacity-achieving decoder from Section A.

## C.1.1. RATE ZERO (R = 0)

Following the form of the decoder analyzed in Section A.1, we consider a test of the form

$$\operatorname{Test}(\hat{\boldsymbol{X}}, \boldsymbol{Y}) = \mathbb{1}\{d^{-1}\langle \boldsymbol{Y}, \hat{\boldsymbol{X}}\rangle \ge 1 - \eta\},\tag{57}$$

where the choice of  $\eta$  will be specified below.

Suppose that  $\hat{\boldsymbol{X}} \in \sqrt{d}\mathcal{S}^{d-1}$  is such that, for some particular  $i \in [k]$ ,  $d^{-1}\|\hat{\boldsymbol{X}} - \boldsymbol{X}_i\|^2 \leq 0.5\varepsilon_{\mathrm{I}}$ . Note that this may be written equivalently as  $d^{-1}\langle \boldsymbol{X}_i, \hat{\boldsymbol{X}} \rangle \geq 1 - 0.25\varepsilon_{\mathrm{I}}$ . Thus,

$$d^{-1}\langle \boldsymbol{X}_i + \sigma \boldsymbol{Z}, \hat{\boldsymbol{X}} \rangle \ge 1 - 0.25\varepsilon_{\mathrm{I}} + \underbrace{(d^{-1/2}\sigma)\langle \boldsymbol{Z}, d^{-1/2}\hat{\boldsymbol{X}}\rangle}_{\sim \mathcal{N}(0, \sigma^2/d)}.$$

Setting

$$\eta = 0.25\varepsilon_{\rm I}\,,\tag{58}$$

we get

$$\begin{split} \Pr(\mathsf{Test}(\hat{\boldsymbol{X}}, \boldsymbol{Y}) = 1) &\geq \frac{1}{k} \Pr(\mathsf{Test}(\hat{\boldsymbol{X}}, \boldsymbol{Y}) = 1 \,|\, \ell = i) \\ &\geq \frac{1}{k} \Pr\left(\mathcal{N}(1 - 0.25\varepsilon_{\mathrm{I}}, \sigma^2/d) \geq 1 - 0.25\varepsilon_{\mathrm{I}}\right) = 0.5/k \,. \end{split}$$

Consequently, with probability 1,  $q_{\text{Close}}(\boldsymbol{\mathcal{X}}_k) \geq 0.5k^{-1}$ .

The challenging part of the analysis is to control  $q_{Far}(\mathcal{X}_k)$ .

Observe that if  $d^{-1}\|\hat{\boldsymbol{X}} - \boldsymbol{X}_i\|^2 \ge \varepsilon_{\mathrm{I}}$  then  $d^{-1}\langle \boldsymbol{X}_i, \hat{\boldsymbol{X}} \rangle \le 1 - 0.5\varepsilon_{\mathrm{I}}$ . For  $\hat{\boldsymbol{X}} \in \sqrt{d}\mathcal{S}^{d-1}$ , denote

$$Q_{i}(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_{k}) = \Pr_{\boldsymbol{Z} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})} \left( \mathsf{Test}(\hat{\boldsymbol{X}}, \boldsymbol{X}_{i} + \sigma \boldsymbol{Z}) = 1 \, \big| \, \boldsymbol{\mathcal{X}}_{k} \right) ,$$

$$\bar{Q}(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_{k}) = \frac{1}{k} \sum_{i=1}^{k} Q_{i}(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_{k}) .$$
(59)

Note that  $Q_i(\hat{X}|\mathcal{X}_k)$  depends on  $\mathcal{X}_k$  only through  $X_i$ .

By definition,  $q_{Far}(\mathcal{X}_k) = \max_{\hat{X} \in \mathcal{T} \cap \mathcal{H}_{Far}} \bar{Q}(\hat{X}|\mathcal{X}_k)$ . We start with a trivial bound.

**Lemma 19** Suppose that  $d^{-1}\langle \mathbf{X}_i, \hat{\mathbf{X}} \rangle \leq 1 - \nu - 0.25\varepsilon_{\mathrm{I}}$  for  $0 \leq \nu \leq 1 - 0.25\varepsilon_{\mathrm{I}}$ . Then

$$Q_i(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_k) \leq k^{-\beta\nu^2}$$
.

Consequently, if  $\hat{X} \in \mathcal{H}_{Far}$  then for all i,  $Q_i(\hat{X}|\mathcal{X}_k) < k^{-\frac{\beta}{16}\varepsilon_1^2}$ .

**Proof** 
$$d^{-1}\langle \boldsymbol{X}_i + \sigma \boldsymbol{Z}, \hat{\boldsymbol{X}} \rangle \leq 1 - 0.25\varepsilon_{\mathrm{I}} - \nu + \underbrace{d^{-1}\sigma\langle \boldsymbol{Z}, \hat{\boldsymbol{X}} \rangle}_{\sim \mathcal{N}(0, \sigma^2/d)}$$
. Thus,

$$\Pr(\mathsf{Test}(\hat{\boldsymbol{X}}, \boldsymbol{X}_i + \sigma \boldsymbol{Z}) = 1) \le \Pr\left(\mathcal{N}(0, \sigma^2/d) \ge \nu\right) \le e^{-\frac{d}{2\sigma^2}\nu^2}.$$

Now, 
$$k=e^{d\mathsf{C}(\beta\sigma^2)}\leq e^{\frac{d}{2\beta\sigma^2}}$$
 (since  $\mathsf{C}(s)=\frac{1}{2}\log(1+1/s)\leq 1/(2s)$ ), therefore  $e^{-\frac{d}{2\sigma^2}\nu^2}\leq k^{-\beta\nu^2}$ . Finally, if  $\hat{\boldsymbol{X}}\in\mathcal{H}_{\mathrm{Far}}$ , then  $d^{-1}\langle\boldsymbol{X}_i,\hat{\boldsymbol{X}}\rangle\leq 1-\nu-0.25\varepsilon_{\mathrm{I}}$  with  $\nu=0.25\varepsilon_{\mathrm{I}}$ .

As mentioned, Lemma 19 gives us the trivial bound  $\bar{Q}(\hat{X}|\mathcal{X}_k) \leq k^{-\frac{\beta}{16}\varepsilon_1^2}$  for all  $\hat{X} \in \mathcal{H}_{Far}$ . This is a highly wasteful bound: it treats  $\hat{X}$  as if it is *simultaneously*  $\sqrt{\varepsilon_1 d}$ -close to all of  $X_1, \ldots, X_k$ . In practice, however, "typical" instances of  $\mathcal{X}_k$  create constellations that do not cluster around any particular point; consequently, for most  $i \in [k]$ , it has to be that, in fact,  $d^{-1}\langle X_i, \hat{X} \rangle \approx 0$ .

For  $t \in (0,1)$ , set

$$N_t(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_k) = \sum_{i=1}^k \mathbb{1}\{d^{-1}\langle \boldsymbol{X}_i, \hat{\boldsymbol{X}}\rangle \ge t\},$$
(60)

the number of centers  $X_i$  that have correlation  $\geq t$  with  $\hat{X}$ .

Choose some constants  $\varepsilon_0, \nu_0 \in (0,1)$  such that  $\nu_0 < 1 - 0.25\varepsilon_0$  and  $\beta\nu_0^2 > 1$ . This can certainly be done, since  $\beta > 1$ . By Lemma 19 above, for every  $\hat{\boldsymbol{X}} \in \mathcal{H}_{Far} \cap \mathcal{T}$ , assuming  $\varepsilon_I \leq \varepsilon_0$ ,

$$\bar{Q}(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_k) \le k^{\beta\nu_0^2} + \frac{1}{k} \cdot N_{(1-\nu_0-0.25\varepsilon_0)}(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_k) \cdot k^{-\frac{\beta}{16}\varepsilon_1^2}.$$
(61)

That is,  $X_i$ -s whose correlation with  $\hat{X}$  is  $<1-\nu_0-0.25\varepsilon_0<1-\nu_0-0.25\varepsilon_{\rm I}$  contribute each at most  $Q_i(\hat{X}|\mathcal{X}_k) \leq k^{-\beta\nu_0^2}$ ; on the other hand, centers whose correlations is higher give, at most, the worst-case contribution  $Q_i(\hat{X}|\mathcal{X}_k) = k^{-\frac{\beta}{16}\varepsilon_1^2}$ . In light of (60), clearly,

$$\mathsf{q}_{\mathrm{Far}}(\boldsymbol{\mathcal{X}}_k) \le k^{\beta\nu_0^2} + \max_{\hat{\boldsymbol{X}} \in \mathcal{T} \cap \mathcal{H}_{\mathrm{Far}}} N_{(1-\nu_0 - 0.25\varepsilon_0)}(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_k) \cdot k^{-1 - \frac{\beta}{16}\varepsilon_1^2} \tag{62}$$

Thus, it remains to show that, with high probability,  $\max_{\hat{X} \in \mathcal{T} \cap \mathcal{H}_{Far}} N_{(1-\nu_0-0.25\varepsilon_0)}(\hat{X}|\mathcal{X}_k)$  is small.

**Lemma 20** Fix any  $\hat{X} \in \mathcal{T}$ . There are universal  $C_1, C_2$  such that whenever  $t \geq C_1 \sqrt{\frac{\log k}{d}}$ , for all  $M \geq 1$ ,

$$\Pr\left(N_t(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_k) \ge M \mid \hat{\boldsymbol{X}} \in \mathcal{H}_{\operatorname{Far}}\right) \le (C_2 k e^{-d\frac{t^2}{2}})^M$$

where the probability is with respect to  $\mathcal{X}_k \sim \mathrm{Unif}(\sqrt{d}\mathcal{S}^{d-1})^{\otimes k}$ , and conditioned on the event that  $\hat{X} \in \mathcal{H}_{\mathrm{Far}}$ .

**Proof** Observe that conditioned on the event  $\hat{X} \in \mathcal{H}_{Far}$ , the centers  $X_1, \dots, X_k$  are i.i.d. and  $\sim \operatorname{Unif}(\sqrt{d}\mathcal{S}^{d-1} \setminus \mathcal{B}(\hat{X}, \sqrt{\varepsilon_1 d}))$ . For any *non-negative*  $f(\cdot)$ ,

$$\mathbb{E}_{\boldsymbol{X}_{i} \sim \text{Unif}(\sqrt{d}\mathcal{S}^{d-1} \setminus \mathcal{B}(\hat{\boldsymbol{X}}, \sqrt{\varepsilon_{I}d}))} [f(\boldsymbol{X}_{i})] \leq \frac{\text{Surf}(\sqrt{d}\mathcal{S}^{d-1})}{\text{Surf}(\sqrt{d}\mathcal{S}^{d-1} \setminus \mathcal{B}(\hat{\boldsymbol{X}}, \sqrt{\varepsilon_{I}d}))} \cdot \mathbb{E}_{\boldsymbol{X}_{i} \sim \text{Unif}(\sqrt{d}\mathcal{S}^{d-1})} [f(\boldsymbol{X}_{i})] \\
= \frac{1}{1 - \varepsilon_{I}^{d-1}} \cdot \mathbb{E}_{\boldsymbol{X}_{i} \sim \text{Unif}(\sqrt{d}\mathcal{S}^{d-1})} [f(\boldsymbol{X}_{i})] \\
= (1 + o(1)) \mathbb{E}_{\boldsymbol{X}_{i} \sim \text{Unif}(\sqrt{d}\mathcal{S}^{d-1})} [f(\boldsymbol{X}_{i})] .$$
(63)

Consequently,  $N_t(\hat{\boldsymbol{X}}, \boldsymbol{\mathcal{X}}_k) \sim \text{Binomial}(n, p)$ , with

$$\begin{split} \mathsf{p} &= \mathbb{E}_{\boldsymbol{X}_i \sim \mathrm{Unif}(\sqrt{d}\mathcal{S}^{d-1} \backslash \mathcal{B}(\hat{\boldsymbol{X}}, \sqrt{\varepsilon_{\mathrm{I}} d}))} \left[ \mathbb{1} \{ d^{-1} \langle \boldsymbol{X}_i, \hat{\boldsymbol{X}} \rangle \geq t \} \right] \\ &= (1 + o(1)) \Pr_{\boldsymbol{X}_i \sim \mathrm{Unif}(\sqrt{d}\mathcal{S}^{d-1})} \left( d^{-1} \langle \boldsymbol{X}_i, \hat{\boldsymbol{X}} \rangle \geq t \right) \\ &\leq 2e^{-d\frac{t^2}{2}} \,, \end{split}$$

where we used the standard tail bound Lemma 35. We have

$$\Pr\left(X_t(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_k) \geq M \mid \hat{\boldsymbol{X}} \in \mathcal{H}_{\operatorname{Far}}\right) = \sum_{l=M}^k \binom{k}{l} \mathsf{p}^l (1-\mathsf{p})^{k-l} \leq \sum_{l=M}^k \left(\frac{ke}{l}\mathsf{p}\right)^l.$$

Assuming  $t \geq C\sqrt{\frac{\log k}{d}}$  for large enough (universal) C > 0,  $\frac{ke}{l} \mathsf{p} \leq \frac{1}{2}$ , and so  $\sum_{l=M}^{k} \left(\frac{ke}{l} \mathsf{p}\right)^{l} \leq 2 \left(ke \mathsf{p}\right)^{M}$ .

**Lemma 21** There are  $C_1, C_2, C_3$  universal such that the following holds.

Suppose that  $t \geq C_1 \sqrt{\frac{\log k}{d}}$ ,  $\varepsilon_{\rm I} < 1/2$ . Then with probability at least  $1 - e^{-d}$  over  $\mathcal{X}_k \sim {\rm Unif}(\sqrt{d}\mathcal{S}^{d-1})^{\otimes k}$ ,

$$\max_{\hat{\boldsymbol{X}} \in \mathcal{T} \cap \mathcal{H}_{\operatorname{Far}}} N_t(\hat{\boldsymbol{X}} | \boldsymbol{\mathcal{X}}_k) \le C_2 \frac{\log(1/\varepsilon_{\mathrm{I}})}{\frac{1}{2}t^2 - \frac{\log k + C_3}{d}} =: M_0.$$

**Proof** Let B be the number of candidates  $\hat{X} \in \mathcal{T}$  such that  $N_t(\hat{X}|\mathcal{X}_k) > M_0$ . Our goal is to show that w.p.  $\geq 1 - e^{-d}$ , B = 0. By Markov's inequality and Lemma 20, assuming  $t \gtrsim \sqrt{\frac{\log k}{d}}$ ,

$$\Pr(B \ge 1) \le \mathbb{E}[B] = (Cke^{-d\frac{t^2}{2}})^{M_0}|\mathcal{T}| \le (Cke^{-d\frac{t^2}{2}})^{M_0}e^{Cd\log(1/\varepsilon_{\rm I})}$$

Taking  $M_0 = C_2 \frac{\log(1/\varepsilon_1)}{\frac{1}{2}t^2 - \frac{\log k + C_3}{d}}$  for large enough  $C_2, C_3$ , the above probability is  $\leq e^{-d}$ .

**Proof** (Of Lemma 9, R = 0.) Recall that, by construction,  $q_{\text{Close}}(\mathcal{X}_k) \geq 0.5k^{-1}$  holds with probability 1.

Recall (62), and apply Lemma 21 with  $t=1-\nu_0-0.25\varepsilon_0$ , which is a positive constant. Consequently, with probability  $\geq 1-e^{-d}$ ,

$$q_{\operatorname{Far}}(\boldsymbol{\mathcal{X}}_k) \leq k^{-\beta\nu_0^2} + C\log(1/\varepsilon_{\operatorname{I}})k^{-1-\frac{\beta}{16}\varepsilon_{\operatorname{I}}^2}$$
.

Since  $\beta \nu_0^2 > 1$ , one may indeed choose some  $c = c(\beta)$  small enough such that  $\mathsf{q}_{\mathrm{Far}}(\mathcal{X}_k) \leq 2k^{-1-c\varepsilon_1^2}$  holds with probability  $1 - o_{\beta,\varepsilon_1}(1)$ .

### C.1.2. Positive Rate (R > 0)

Moving on to the positive rate regime, our construction is guided by the decoder of Section A.2.

Set 
$$\alpha = \frac{1}{1+\sigma^2}$$
,  $\tau = \frac{\sigma^2}{1+\sigma^2}$ .

For fixed  $\hat{X}, X_i \in \sqrt{d}S^{d-1}$ , denoting  $Y = X_i + \sigma Z, Z \sim \mathcal{N}(0, I)$ ,

$$\mathbb{E}_{\boldsymbol{Z}} \|\alpha \boldsymbol{Y} - \hat{\boldsymbol{X}}\|^{2} = \mathbb{E}_{\boldsymbol{Z}} \|\alpha \boldsymbol{Y} - \boldsymbol{X}_{i}\|^{2} + \|\boldsymbol{X}_{i} - \hat{\boldsymbol{X}}\|^{2} + 2\mathbb{E}_{\boldsymbol{Z}} \langle \alpha \boldsymbol{Y} - \boldsymbol{X}_{i}, \boldsymbol{X}_{i} - \hat{\boldsymbol{X}} \rangle$$

$$= \tau d + \|\boldsymbol{X}_{i} - \hat{\boldsymbol{X}}\|^{2} + 2(\alpha - 1)\langle \boldsymbol{X}_{i}, \boldsymbol{X}_{i} - \hat{\boldsymbol{X}} \rangle$$

$$= \tau d + \alpha \|\boldsymbol{X}_{i} - \hat{\boldsymbol{X}}\|^{2},$$
(64)

where we used  $\|\boldsymbol{X}_i\|^2 = \|\hat{\boldsymbol{X}}\|^2 = d$  and  $-\langle \boldsymbol{X}_i, \hat{\boldsymbol{X}} \rangle = \frac{1}{2}(\|\boldsymbol{X}_i - \hat{\boldsymbol{X}}\|^2 - \|\boldsymbol{X}_i\|^2 - \|\hat{\boldsymbol{X}}\|^2)$ .

Assume that  $\|\boldsymbol{X}_i - \hat{\boldsymbol{X}}\|^2 \leq 0.5\varepsilon_{\mathrm{I}}d$ . By the Gaussian Lipschitz concentration inequality (Lemma 34), applied for  $F(\boldsymbol{Z}) = d^{-1/2}\|\alpha(\boldsymbol{X}_i + \sigma\boldsymbol{Z}) - \hat{\boldsymbol{X}}\|$ , which is  $(d^{-1/2}\alpha\sigma)$ -Lipschitz with expectation  $\mathbb{E}_{\boldsymbol{Z}}F(\boldsymbol{Z}) \leq \sqrt{\mathbb{E}(F(\boldsymbol{Z}))^2} \leq \sqrt{\tau + 0.5\alpha\varepsilon_0}$ ,

$$\Pr_{\mathbf{Z}}\left(d^{-1/2}\|\alpha\mathbf{Y}_i - \hat{\mathbf{X}}\| \ge \sqrt{\tau + 0.5\alpha\varepsilon_{\mathrm{I}}} + \eta\right) \le e^{-\frac{1}{2}\frac{\eta^2}{\alpha^2\sigma^2}d}.$$

Consider the test

$$\operatorname{Test}(\hat{\boldsymbol{X}}, \boldsymbol{Y}) = \mathbb{1}\{d^{-1/2}\|\alpha \boldsymbol{Y} - \hat{\boldsymbol{X}}\| \le \sqrt{\tau + 0.5\alpha\varepsilon_I} + \eta\}, \quad \eta = \sqrt{\frac{2\alpha^2\sigma^2\log 2}{d}} = O_{\mathsf{R},\beta}(d^{-1/2}),$$
(65)

so that by construction,  $q_{\text{Close}}(\mathcal{X}_k) \geq 0.5k^{-1}$  holds with probability 1.

It remains to bound  $q_{Far}(\mathcal{X}_k)$  with high probability. We follow the notation (59), where  $\mathsf{Test}(\hat{X}, Y)$  that appears in (59) is now defined by (65). Our goal is to bound, with high probability over  $\mathcal{X}_k$ ,

$$\mathsf{q}_{\mathrm{Far}}(\boldsymbol{\mathcal{X}}_k) = \max_{\hat{\boldsymbol{X}} \in \mathcal{T} \cap \mathcal{H}_{\mathrm{Far}}} \bar{Q}_i(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_k) = \max_{\hat{\boldsymbol{X}} \in \mathcal{T} \cap \mathcal{H}_{\mathrm{Far}}} \frac{1}{k} \sum_{i=1}^k Q_i(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_k) \,.$$

As was in the zero rate case, for any fixed  $\hat{X}$ , conditioned on the event  $\{\hat{X} \in \mathcal{H}_{\mathrm{Far}}\}$ , the centers  $X_1, \ldots, X_k$  are i.i.d. and  $X_i \sim \mathrm{Unif}(\sqrt{d}\mathcal{S}^{d-1} \setminus \mathcal{B}(\hat{X}, \sqrt{\varepsilon_{\mathrm{I}}d}))$ . Consequently,  $\bar{Q}(\hat{X}|\mathcal{X}_k)$  is the average of k i.i.d. random variables. We shall show that its expectation is very small, specifically  $\mathbb{E}[\bar{Q}(\hat{X}|\mathcal{X}_k)] \leq k^{-1-c}$ ; moreover, we shall show that it concentrates tightly about this expectation, to the extent that the maximum over the net,  $\max_{\hat{X} \in \mathcal{T} \cap \mathcal{H}_{\mathrm{Far}}} \bar{Q}_i(\hat{X}|\mathcal{X}_k)$ , can be controlled as well. To do this, we use Bernstein's inequality for sums of i.i.d. bounded random variables, Lemma 36.

For brevity, we introduce some notation. For  $\hat{X} \in \mathcal{T}$  fixed, denote by  $\mathcal{E} = \mathcal{E}(\hat{X})$  the event  $\mathcal{E} = \{\hat{X} \in \mathcal{H}_{\operatorname{Far}}\}$  (with respect to the probability on  $\mathcal{X}_k \sim \operatorname{Unif}(\sqrt{d}\mathcal{S}^{d-1})^{\otimes k}$ ). Denote by  $\mathbb{E}^{\mathcal{E}}[\cdot], \|\cdot\|_{\infty}^{\mathcal{E}}$  respectively the expectation and  $L_{\infty}$  norm with respect to the conditional measure on  $\mathcal{X}_k$ ; and  $\mathbb{P}^{\mathcal{E}}(S) := \mathbb{E}^{\mathcal{E}}[\mathbb{1}_S]$ .

To use Bernstein's inequality, we need two components: an  $L_{\infty}$  bound and a bound on the expectation.

We start with the  $L_{\infty}$  bound:

**Lemma 22** There are  $C, \varepsilon_0$ , that depend on  $R, \beta$ , such that the following holds. For any  $\hat{X} \in \mathcal{T}$ , whenever  $\varepsilon_I < \varepsilon_0$  and d is sufficiently large,  $\varepsilon_I = \Omega_{R,\beta}(d^{-1/2})$ , then

$$\|Q_i(\hat{\mathbf{X}}|\mathcal{X}_k)\|_{\infty}^{\mathcal{E}} \le 2k^{-C\varepsilon_I^2}$$
 (66)

**Proof** Fix any  $\hat{X}$ ,  $\mathcal{X}_k$  such that  $\hat{X} \in \mathcal{H}_{Far}$ . Let  $\mu = d^{-1/2} \|\alpha X_i - \hat{X}\| \leq 2$ . By the rotational invariance of  $Z \sim \mathcal{N}(\mathbf{0}, I)$ ,

$$\|\alpha(\boldsymbol{X}_i + \sigma \boldsymbol{Z}) - \hat{\boldsymbol{X}}\|^2 \stackrel{d}{=} \|\alpha \sigma \boldsymbol{Z} + \mu \boldsymbol{1}\|^2 = \sum_{j=1}^d (\alpha \sigma Z_j + \mu)^2,$$

where  $\mathbf{1}=(1,\ldots,1)$  and  $\stackrel{d}{=}$  denotes equality in distribution. The expression above is a sum of i.i.d. sub-Exponential random variables. The sub-Exponential norm, denoted  $\|\cdot\|_{\psi_1}$ , is upper bounded by

$$\|(\alpha\sigma Z_j + \mu)^2\|_{\psi_1} \lesssim \alpha^2 \sigma^2 + \mu^2 = O_{\beta,\mathsf{R}}(1).$$

For background on sub-Exponential random variables, including the definition of the sub-Exponential norm (and Orlicz norms in general), we refer to the book (Vershynin, 2018, Chapter 2). Recall by (64) that

$$d^{-1}\mathbb{E}\|\alpha\sigma\mathbf{Z} + \mu\mathbf{1}\|^2 = \tau + \alpha d^{-1}\|\mathbf{X}_i - \hat{\mathbf{X}}\|^2 \ge \tau + \alpha\varepsilon_{\mathrm{I}},$$

and set

$$t = \tau + \alpha \varepsilon_{\rm I} - \left(\sqrt{\tau + 0.5\alpha\varepsilon_{\rm I}} + \eta\right)^{2}$$
$$= \tau + \alpha \varepsilon_{\rm I} - \left(\sqrt{\tau + 0.5\alpha\varepsilon_{\rm I}} + O_{\rm R,\beta}(d^{-1/2})\right)^{2}$$
$$= 0.5\alpha \varepsilon_{\rm I} + O_{\rm R,\beta}(d^{-1/2}).$$

Using Bernstein's inequality for sub-Exponential random variables, Lemma 37,

$$Q_{i}(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_{k}) = \Pr_{\boldsymbol{Z}}(\mathsf{Test}(\hat{\boldsymbol{X}}, \boldsymbol{X}_{i} + \sigma \boldsymbol{Z}) = 1)$$

$$= \Pr\left(d^{-1}\|\alpha(\boldsymbol{X}_{i} + \sigma \boldsymbol{Z}) - \hat{\boldsymbol{X}}\|^{2} \leq \left(\sqrt{\tau + 0.5\alpha\varepsilon_{I}} + \eta\right)^{2}\right)$$

$$\leq \Pr\left(d^{-1}\|\alpha(\boldsymbol{X}_{i} + \sigma \boldsymbol{Z}) - \hat{\boldsymbol{X}}\|^{2} - \mathbb{E}_{\boldsymbol{Z}}d^{-1}\mathbb{E}_{\boldsymbol{Z}}\|\alpha(\boldsymbol{X}_{i} + \sigma \boldsymbol{Z}) - \hat{\boldsymbol{X}}\|^{2} \leq -t\right)$$

$$\leq 2\exp\left(-c\min\left\{\frac{(d \cdot t)^{2}}{d \cdot \|\alpha\sigma\boldsymbol{Z}_{j} + \mu\|_{\psi_{1}}^{2}}, \frac{d \cdot t}{\|\alpha\sigma\boldsymbol{Z}_{j} + \mu\|_{\psi_{1}}}\right\}\right)$$

$$\stackrel{(\star)}{\leq} 2\exp\left(-C(\mathsf{R}, \beta)\min\{\varepsilon_{I}, \varepsilon_{I}^{2}\} \cdot d\right)$$

$$= 2\exp\left(-C(\mathsf{R}, \beta)\varepsilon_{I}^{2} \cdot d\right)$$

where to get  $(\star)$ , we used  $\|\alpha\sigma Z_j + \mu\|_{\psi_1} = O_{\beta,R}(1)$  and  $\varepsilon_1 \gtrsim d^{-1/2}$ . Finally, to deduce (66), recall that  $k = e^{R \cdot d}$ .

Moving on to the expectation:

**Lemma 23** For any  $\varepsilon_{\rm I} < \varepsilon_0$  sufficiently small and d sufficiently large such that  $\varepsilon_{\rm I} = \Omega_{{\sf R},\beta}(d^{-1/2})$ ,

$$\mathbb{E}^{\mathcal{E}}[Q_i(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_k)] \le 2k^{-1-c},$$

where  $c, \varepsilon_0$  depend on  $R, \beta$ .

<sup>3.</sup> An alternative method to the one below (which is itself very standard) is to use deviation inequalities for non-central  $\chi^2$  random variables, that are readily available in the literature, though somewhat more "messy".

**Proof** Observe that the test, defined in (65), is an orthogonally invariant function of its argument; that is,  $\mathsf{Test}(\boldsymbol{x}, \boldsymbol{y}) = \mathsf{Test}(R\boldsymbol{x}, R\boldsymbol{y})$  for any  $R \in O(d)$ , where O(d) is the group of  $d \times d$  orthogonal matrices. Introduce an independent  $R \sim \mathsf{Haar}(O(d))$ , and note that, by the orthogonal invariance of  $\boldsymbol{Z} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})$ ,  $R(\boldsymbol{X}_i + \sigma \boldsymbol{Z}) \stackrel{d}{=} R\boldsymbol{X}_i + \sigma \boldsymbol{Z}$ .

Now, conditioned on  $RX_i$ , the conditional distribution of  $R\hat{X}$  is  $\sim \text{Unif}(\sqrt{d}S^{d-1}\setminus\mathcal{B}(RX_i,\sqrt{\varepsilon_Id}))$ . Thus, again owing to orthogonal invariance, we have

$$\mathbb{E}^{\mathcal{E}}[Q_{i}(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_{k})] = \Pr_{\boldsymbol{W} \sim \text{Unif}(\sqrt{d}S^{d-1} \setminus \mathcal{B}(\boldsymbol{x}, \sqrt{\varepsilon_{I}d})), \\ \boldsymbol{Z} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})} \left( \|\alpha(\boldsymbol{x} + \sigma \boldsymbol{Z}) - \boldsymbol{W}\| \leq \sqrt{\tau + 0.5\varepsilon_{I}} + \eta \right) \\ \stackrel{(\star)}{\leq} \frac{1}{1 - \varepsilon_{I}^{d-1}} \Pr_{\boldsymbol{W} \sim \text{Unif}(\sqrt{d}S^{d-1}), \\ \boldsymbol{Z} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})} \Pr\left( \|\alpha(\boldsymbol{x} + \sigma \boldsymbol{Z}) - \boldsymbol{W}\| \leq \sqrt{\tau + 0.5\varepsilon_{I}} + \eta \right) ,$$

where  $x \in \sqrt{d}S^{d-1}$  is any fixed vector, and  $(\star)$  follows from (63). The probability above has been bounded in Lemma 17, which yields

$$\leq (1+o(1))\left(1+\frac{1}{\sigma^2}\right)^{1/2}e^{-\left(\mathsf{C}(\sigma^2)-\frac{1}{2}\log\frac{(\sqrt{\tau+0.5\varepsilon_I}+\eta)^2}{\tau}\right)d}\,.$$

Since  $k = e^{\mathsf{C}(\beta\sigma^2)d}$  for  $\beta > 1$  constant (hence  $\mathsf{C}(\sigma^2) - \mathsf{C}(\beta\sigma^2)$  is a positive constant), and assuming  $\varepsilon_{\mathrm{I}} \leq \varepsilon_0$  is small enough, we get that  $\mathbb{E}^{\mathcal{E}}[Q_i(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_k)] \leq 2k^{-1-c}$  for some  $c = c(\beta,\mathsf{R})$ , for d large enough.

We are ready to bound  $q_{Far}(\boldsymbol{\mathcal{X}}_k)$ :

**Lemma 24** There are  $C, \varepsilon_0$  that depend on  $R, \beta$ , such that whenever  $\varepsilon_1 < \varepsilon_0$  then

$$\mathsf{q}_{\mathrm{Far}}(\boldsymbol{\mathcal{X}}_k) \leq 3k^{-1-C\varepsilon_{\mathrm{I}}^2}$$

holds with probability  $1 - o_{\varepsilon_{\mathbf{I}}, \mathsf{R}, \beta}(1)$  over  $\mathcal{X}_k \sim \mathrm{Unif}(\sqrt{d}\mathcal{S}^{d-1})^{\otimes k}$ .

**Proof** Fix  $\hat{X} \in \mathcal{T}$ . We start by showing that conditioned on  $\hat{X} \in \mathcal{H}_{Far}$ ,  $\bar{Q}(\hat{X}|\mathcal{X}_k)$  is very small with high probability; to that end, we shall use Bernstein's inequality, Lemma 36.

By Lemma 22,  $\|Q_i(\boldsymbol{\mathcal{X}}_k|\hat{\boldsymbol{X}})\|_{\infty}^{\mathcal{E}} \leq 2k^{-C_1\varepsilon_I^2}$ . By Lemma 23,  $\mathbb{E}^{\mathcal{E}}[Q_i(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_k)] \leq 2k^{-1-C_2}$ . Consequently, for small enough  $\varepsilon_I$ ,  $\|Q_i(\boldsymbol{\mathcal{X}}_k|\hat{\boldsymbol{X}}) - \mathbb{E}^{\mathcal{E}}[Q_i(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_k)]\|_{\infty}^{\mathcal{E}} \leq 4k^{-C_1\varepsilon_I^2}$ . Note moroever that since  $Q_i(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_k) \geq 0$ ,

$$\operatorname{Var}^{\mathcal{E}}(Q_{i}(\boldsymbol{\mathcal{X}}_{k}|\hat{\boldsymbol{X}})) \leq \mathbb{E}^{\mathcal{E}}((Q_{i}(\boldsymbol{\mathcal{X}}_{k}|\hat{\boldsymbol{X}})^{2}) \leq \|Q_{i}(\boldsymbol{\mathcal{X}}_{k}|\hat{\boldsymbol{X}})\|_{\infty}^{\mathcal{E}}\mathbb{E}^{\mathcal{E}}[Q_{i}(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_{k})] \leq 4k^{-1-C_{2}-C_{1}\varepsilon_{1}^{2}}.$$

Since  $Q_1(\hat{X}|\mathcal{X}_k), \dots, Q_k(\hat{X}|\mathcal{X}_k)$  are i.i.d. conditioned on  $\mathcal{E}$ , by Bernstein's inequality for sums of independent bounded random variables, for some universal c,

$$\mathbb{P}^{\mathcal{E}}(\bar{Q}(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_k) \geq t + 2k^{-1-C_2}) \leq 2\exp\left(-ck\min\left\{\frac{t^2}{\mathrm{Var}^{\mathcal{E}}(Q_i(\boldsymbol{\mathcal{X}}_k|\hat{\boldsymbol{X}}))}, \frac{t}{\|Q_i(\boldsymbol{\mathcal{X}}_k|\hat{\boldsymbol{X}}) - \mathbb{E}^{\mathcal{E}}[Q_i(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_k)]\|_{\infty}^{\mathcal{E}}}\right\}\right).$$

Setting  $t=k^{1-C_1\varepsilon_1^2/2}$  and assuming  $\varepsilon_I$  is small enough, for c (perhaps other) universal,

$$\mathbb{P}^{\mathcal{E}}(\bar{Q}(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_k) \ge 3k^{-1-C_1\varepsilon_1^2/2}) \le \exp\left(-ck^{C_1\varepsilon_1^2/2}\right).$$

So far we have shown that with overwhelming probability over the configuration  $\mathcal{X}_k$ , conditioned on  $\hat{X} \in \mathcal{H}_{\operatorname{Far}}$ ,  $\bar{Q}(\hat{X}|\mathcal{X}_k)$  is very small. We now wish to control  $\operatorname{q_{\operatorname{Far}}}(\mathcal{X}_k) = \max_{\hat{X} \in \mathcal{T} \cap \mathcal{H}_{\operatorname{Far}}} \bar{Q}(\hat{X}|\mathcal{X}_k)$ . Let  $N = \sum_{\hat{X} \in \mathcal{T}} \mathbb{1}\{\hat{X} \in \mathcal{H}_{\operatorname{Far}} \text{ and } \bar{Q}(\hat{X}|\mathcal{X}_k) \geq 3k^{-1-C_1\varepsilon_1^2/2}\}$ . Of course, N = 0 implies that  $\operatorname{q_{\operatorname{Far}}}(\mathcal{X}_k) \leq 3k^{-1-C_1\varepsilon_1^2/2}$ . By Markov's inequality,

$$\Pr(N \geq 1) \leq \mathbb{E}[N] \leq |\mathcal{T}| \mathbb{P}^{\mathcal{E}}(\bar{Q}(\hat{\boldsymbol{X}}|\boldsymbol{\mathcal{X}}_k) \geq 3k^{-1-C_1\varepsilon_1^2/2}) \leq e^{Cd\log(1/\varepsilon_1)} \exp\left(-ck^{C_1\varepsilon_1^2/2}\right) \; .$$

Recall that k is exponential in d; consequently, for  $\varepsilon_{\rm I} \geq C_2 \sqrt{\frac{\log d}{d}}$  with large enough  $C_2$  (in particular, whenever  $\varepsilon_{\rm I} > 0$  is constant),

$$e^{Cd\log(1/\varepsilon_{\mathrm{I}})}\exp\left(-ck^{C_1\varepsilon_{\mathrm{I}}^2/2}\right) \leq e^{-Cd\log d}\exp(-cd^2) = o_{\varepsilon_{\mathrm{I}},\beta,\mathsf{R}}(1)\,.$$

**Proof** (Of Lemma 9). Follows immediately from Lemma 24, where note that since  $q_{Far}(\mathcal{X}_k) \leq 1$ , we may change the prefactor 3 in Lemma 24 to whatever number > 1 we like (at the expense of changing the exponent). We do so for convenience.

### C.2. Proof of Lemma 10

We start by showing that w.h.p.,  $\mathcal{T}_{\text{Close}} \cap \mathcal{H}_{\text{Far}} = \emptyset$ , namely, we do not retain candidates which are  $\sqrt{\varepsilon_1 d}$ -far from all centers  $X_i$ .

Let  $\hat{X} \in \mathcal{T} \cap \mathcal{H}_{\mathrm{Far}}$ . Conditioned on the event of Lemma 9, the random variables  $\{\mathrm{Test}(\hat{X}, Y)\}$  are an i.i.d. sequence of N Bernoulli trials, with success probability  $\leq 2k^{-1-c\varepsilon_1^2}$ . By Chernoff's inequality (Lemma 38) and the estimate of Lemma 39,

$$\Pr(\hat{\boldsymbol{X}} \in \mathcal{T}_{\text{Close}}) \le e^{-C_1 N \varepsilon_1^2 k^{-1} \log k}$$

for some  $C_1 = C_1(R, \beta)$ , and assuming k is large enough. Taking a union bound,

$$\Pr(\mathcal{T}_{\text{Close}} \cap \mathcal{H}_{\text{Far}} \neq \emptyset) \leq |\mathcal{H}_{\text{Far}} \cap \mathcal{T}| e^{-C_1 N \varepsilon_{\text{I}}^2 k^{-1} \log k} \leq e^{Cd \log(1/\varepsilon_{\text{I}})} e^{-C_1 N \varepsilon_{\text{I}}^2 k^{-1} \log k}.$$

Observe that this is  $o_{R,\beta}(1)$  whenever  $N \geq C_2 \frac{\log(1/\varepsilon_1)}{\varepsilon_1^2} k \frac{d}{\log k}$ , for some  $C_2 = C_2(R,\beta)$  large enough. Now, when R > 0 then  $\frac{d}{\log k} = 1/R_{d,k} = 1/R$  is just a constant, and so is  $\sigma^2$ ; so for a suitably modified  $C_3 = C_3(R,\beta)$ ,  $N \geq C_3 \frac{\log(1/\varepsilon_1)}{\varepsilon_1^2} k \sigma^2$  suffices. As for the case R = 0,  $R_{d,k} = C(\beta\sigma^2) = \frac{1}{\beta\sigma^2} + O(\sigma^{-4})$ , and so  $\frac{d}{\log k} \lesssim \sigma^2$ ; consequently,  $N \geq C_3 \frac{\log(1/\varepsilon_1)}{\varepsilon_1^2} k \sigma^2$  suffices. This show the first claim of the Lemma.

Moving on, we need to show that for many  $i \in [k]$ ,  $\mathcal{T}_{\text{Close}}$  indeed contains a vector within  $\sqrt{\varepsilon_{\text{I}}d}$ -distance to  $X_i$ .

Fix any  $\hat{X}_1, \ldots, \hat{X}_k \in \mathcal{T}$  such that  $\|X_i - \hat{X}_i\|^2 \leq 0.5\varepsilon_{\mathrm{I}}d$ . Since  $\mathcal{T}$  is an  $\sqrt{0.5\varepsilon_{\mathrm{I}}d}$ -net of  $\sqrt{d}\mathcal{S}^{d-1}$ , there certainly are such vectors in  $\mathcal{T}$ . We shall show that  $|\mathcal{T}_{\mathrm{Close}} \cap \{\hat{X}_1, \ldots, \hat{X}_k\}| \geq (1-\varphi)k$  holds with the claimed probability; this clearly suffices. By the properties of the test, Lemma 9,  $\Pr(\mathsf{Test}(\hat{X}_i, Y) = 1) \geq 0.5k^{-1}$ . By Chernoff's inequality (Lemma 38),

$$\Pr(\hat{X}_i \notin \mathcal{T}_{\text{Close}}) \le e^{-cNk^{-1}}$$
.

for some universal c. Consequently, by Markov's inequality,

$$\Pr\left(\sum_{i=1}^{k} \mathbb{1}\{\hat{\boldsymbol{X}}_i \notin \mathcal{T}_{\text{Close}}\} \ge \varphi k\right) \le \varphi^{-1} e^{-cNk^{-1}}.$$

Thus, when  $N \ge Ck \log(1/\varphi)$ , for some universal C, we get that with probability  $\ge 1 - \varphi$ , we have  $|\mathcal{T}_{\text{Close}} \cap \{\hat{X}_1, \dots, \hat{X}_k\}| \ge (1 - \varphi)k$ .

## C.3. Proof of Lemma 11

As discussed in the main text, it suffices to show that with probability  $1 - o_{\beta,R}(1)$ , the centers in  $\mathcal{X}_k$  have minimal distance  $\sqrt{Ld}$  for some  $L = L(\beta,R)$ . In that case, choosing  $\varepsilon_0 \leq L/16$ , the required results follows immediately from Lemmas 9 and 10.

The following argument is standard. Sample centers  $X_1,\ldots,X_k$  sequentially. Let  $E_i$  be the event that  $X_i \notin \bigcup_{l=1}^{i-1} \mathcal{B}(X_l,\sqrt{Ld})$  for  $i=2,\ldots,k$ . Clearly,  $\mathcal{X}_k$  has minimal distance  $\geq \sqrt{Ld}$  if and only if  $\bigcap_{l=2}^k E_i$  holds. Notice that  $\Pr(E_i|\bigcap_{l=1}^{i-1} E_i) \geq 1 - L^{d-1}(i-1)$ , since  $X_i$  has to evade i-1 disjoint neighborhoods  $\sqrt{d}\mathcal{S}^{d-1} \cap \mathcal{B}(X_i,\sqrt{RLd})$ , that amount to total surface area at most  $\leq (i-1) \mathrm{Surf}(\partial \mathcal{B}(X_i,\sqrt{Ld})) = (i-1) L^{d-1} \mathrm{Surf}(\sqrt{d}\mathcal{S}^{d-1})$ . Somewhat crudely, we lower bound:

$$\Pr\left(\bigcap_{i=2}^{k} E_{i}\right) = \prod_{i=2}^{k} \Pr(E_{l} | \bigcap_{l=2}^{i-1} E_{l}) \ge (1 - L^{d-1}k)^{k}.$$
(67)

Whenever  $L^{d-1}k = o(1/k)$ , the bound in (67) tends to 1 as  $k \to \infty$ . When  $k = e^{o(d)}$ , any constant L < 1 will work. When  $k = e^{dR}$ , any constant  $L < e^{-2R}$  will work.

## **Appendix D. Proofs for Section 4.2**

### D.1. Decoding Using a Corrupted Codebook

Upon successful completion of Step I, by Lemma 11, we will have constructed a list  $\tilde{\mathcal{X}}_{\mathrm{I}} = (\tilde{X}_1, \dots, \tilde{X}_m)$  of size  $m \geq (1-\varphi)k$ , such that every member of  $\tilde{\mathcal{X}}_{\mathrm{I}}$  is  $\sqrt{\varepsilon_{\mathrm{I}}d}$ -close to some unique codeword of  $\mathcal{X}_k$ . In this section, we show that whenever  $\varepsilon_{\mathrm{I}}$  is smaller than *some particular threshold*  $\varepsilon_0 = \varepsilon_0(\beta, \mathbb{R})$ , then  $\tilde{\mathcal{X}}_{\mathrm{I}}$  may be used to successfully decode messages encoded with  $\mathcal{X}_k$ , in the following sense:

- Whenever  $\tilde{\mathcal{X}}_{I}$  contains a point  $\tilde{X}_{l}$  close to  $X_{i(l)} \in \mathcal{X}_{k}$ , applying the decoder on observation  $Y = X_{i(l)} + \sigma Z$  will indeed return, with high probability, the correct index l.
- As importantly, whenever  $X_I$  does not contain a keyword which is close to  $X_i$ , then applying the decoder on observation  $Y = X_i + \sigma Z$  will consistently return an error symbol "#"; that is, the decoder will not erroneously assign a sample Y to a label even if it does not have a close approximation for its corresponding center  $X_i$ .

We now proceed to formalize the discussion above.

Denote by

$$\operatorname{Approx}(\mathcal{X}_k) \subseteq \bigcup_{m=0}^k (\sqrt{d}\mathcal{S}^{d-1})^{\otimes m}$$
(68)

the set of all lists  $\tilde{\mathbf{X}} = (\tilde{\mathbf{X}}_1, \dots, \tilde{\mathbf{X}}_m)$  (for any  $0 \le m \le k$ ), for which there exists a permutation  $\mathbf{i} : [k] \to [k]$  satisfying

$$\|\tilde{\boldsymbol{X}}_l - \boldsymbol{X}_{\mathsf{i}(l)}\|^2 \le \varepsilon d \quad \text{for all } 1 \le l \le m.$$
 (69)

A family of decoders is a mapping  $\tilde{\boldsymbol{\mathcal{X}}} \mapsto \mathsf{Dec}(\cdot|\tilde{\boldsymbol{\mathcal{X}}})$ , mapping codebooks of any length  $0 \le m \le k$  to decision rules  $\mathbb{R}^d \to [k] \cup \{\#\}$ .  $\mathsf{Dec}(\cdot|\tilde{\boldsymbol{\mathcal{X}}})$  may depend on  $d, k, \sigma^2$  as well, and this shall be implied from now on.

For example, consider the nearest-neighbor family of decoders:

$$DecNN(Y|\tilde{X}) = \underset{l \in [m]}{\operatorname{argmin}} \|Y - \tilde{X}_l\|.$$
 (70)

Recall that when  $\tilde{\mathcal{X}} = \mathcal{X}_k$ , the resulting decoder  $\operatorname{DecNN}(\cdot|\mathcal{X}_k)$  is optimal (in the sense of average error) for decoding a message  $\ell \sim \operatorname{Unif}([k])$  encoded using  $\mathcal{X}_k$ . However, when  $\tilde{\mathcal{X}}$  is only a partial sub-codebook of  $\mathcal{X}_k$ , using the decoder (70) might not be a good idea from a practical standing: an observation  $Y = X_i + \sigma Z$  corresponding to a codeword  $X_i$  which is absent from  $\tilde{\mathcal{X}}$  will necessarily be decoded into an erroneous message. It is desirable that having identifed such a case, the decoder would instead declare an error.

For a codebook  $\mathcal{X}_k \in (\sqrt{d}\mathcal{S}^{d-1})^{\otimes k}$  and  $\tilde{\mathcal{X}} \in \operatorname{Approx}_{\varepsilon}(\mathcal{X}_k)$ , let  $i:[k] \to [k]$  be the permutation that satisfies (69). When there is more than one such permutation (as will surely be the case when  $m \leq k-2$ ), suppose that i is chosen in some systematic way, such that the assignment  $\operatorname{Approx}_{\varepsilon}(\mathcal{X}_k) \to \operatorname{Sym}(k)$ ,  $\tilde{\mathcal{X}} \mapsto i$  is well-defined. Note that  $\mathrm{i}([m])$  are the indices  $\subseteq [k]$  of codewords  $X_i$  for which  $\tilde{\mathcal{X}}$  contains an  $\sqrt{\varepsilon d}$ -distance approximation.

Similar to (17), we consider the error probability of decoding a message  $i \in [k]$ , encoded using codebook  $\mathcal{X}_k$ , with a decoder  $\mathsf{Dec}(\cdot|\tilde{\mathcal{X}})$ ,  $\tilde{\mathcal{X}} \in \mathsf{Approx}_{\varepsilon}(\mathcal{X}_k)$ :

$$P_{approx,i}(\sigma^{2}|\boldsymbol{\mathcal{X}}_{k},\tilde{\boldsymbol{\mathcal{X}}},\mathsf{Dec}(\cdot|\cdot)) = \begin{cases} \Pr(\mathsf{i}^{-1}(i) \neq \mathsf{Dec}(\boldsymbol{X}_{i} + \sigma \boldsymbol{Z}|\tilde{\boldsymbol{\mathcal{X}}})) & \text{if } i \in \mathsf{i}([m]) \\ \Pr(\# \neq \mathsf{Dec}(\boldsymbol{X}_{i} + \sigma \boldsymbol{Z}|\tilde{\boldsymbol{\mathcal{X}}})) & \text{if } i \notin \mathsf{i}([m]) \end{cases} . \tag{71}$$

The "twist" over (17) is that if  $i \notin i([m])$ , we consider decoding to be successful if the decoder declares error.

We are ready to state the technical result of this section.

**Proposition 25** Suppose that  $\beta > 1$ ,  $R_{d,k} = C(\beta \sigma^2)$  (in either positive or zero rate, as before),  $\mathcal{X}_k \sim \mathrm{Unif}(\sqrt{d}\mathcal{S}^{d-1})^{\otimes k}$ .

*There is*  $\varepsilon_0 = \varepsilon_0(\beta, R)$  *and a family of decoders*  $Dec(\cdot|\cdot)$  *such that for all*  $\varepsilon \leq \varepsilon_0$ :

$$\lim_{d\to\infty} \mathbb{E}\left[\sup_{\tilde{\boldsymbol{X}}\in \mathrm{Approx}_{\varepsilon}(\boldsymbol{\mathcal{X}}_k)} \frac{1}{k} \sum_{i=1}^k P_{approx,i}(\sigma^2|\boldsymbol{\mathcal{X}}_k,\tilde{\boldsymbol{\mathcal{X}}}, \mathrm{Dec}(\cdot|\cdot))\right] = 0\,.$$

In words: Proposition 25 states that provided that  $\varepsilon$  deceeds some particular threshold, one can construct a decoder family  $\operatorname{Dec}(\cdot|\cdot)$ , so that for any *adversarially chosen*  $\tilde{\mathcal{X}} \in \operatorname{Approx}_{\varepsilon}(\mathcal{X}_k)$ , the average decoding error (in the sense of (71)) is uniformly small; and that this holds for "most" random codebooks  $\mathcal{X}_k$ .

### D.2. Proof of Proposition 25

We simply adapt the decoders used in the proof of Proposition 4, and appearing in Section A. We give a different construction at zero (R = 0) and positive (R > 0) rate.

## D.2.1. ZERO RATE (R = 0)

We adapt the decoder from Section A.1. Recall the decision rule implemented by this decoder (assuming access to the true codebook  $\mathcal{X}_k$ ): it returns  $i \in [k]$  if and only if: 1)  $d^{-1}\langle \mathbf{Y}, \mathbf{X}_i \rangle \geq 1 - \eta_1$ ; 2) For all  $j \in [k] \setminus \{i\}$ ,  $d^{-1}\langle \mathbf{Y}, \mathbf{X}_j \rangle \leq 1 - \eta_2$ ; if no such i exists it returns an error. Here,  $0 < \eta_1 < \eta_2$  are sufficiently small constants.

From now on, suppose without loss of generality that a message i is sent. Denote by  $Y = X_i + \sigma Z$ , the channel output. Let  $S_i$  the event that the above decoder (which utilizes  $\mathcal{X}_k$ ) succeeds; in other words,

$$S_i(\eta_1, \eta_2) = \{ d^{-1} \langle \boldsymbol{Y}, \boldsymbol{X}_i \rangle \ge 1 - \eta_1 \} \cap \bigcap_{j \in [k] \setminus \{i\}} \{ d^{-1} \langle \boldsymbol{Y}, \boldsymbol{X}_i \rangle \le 1 - \eta_2 \}.$$
 (72)

Recall: in Section A.1 we proved that  $\mathbb{E}_{Y_i, \mathcal{X}_k}[\mathbb{1}_{S_i(\eta_1, \eta_2)}] = 1 - o(1)$  for all sufficiently small constants  $0 < \eta_1 < \eta_2$ .

We now adapt the construction of Section A.1 to use the "corrupted" codebook  $\mathcal{X}$ . For thresholds  $\tilde{\eta}_1, \tilde{\eta}_2$ , the decoder returns l if and only if  $\tilde{X}_l$  is such that  $d^{-1}\langle Y, \tilde{X}_l \rangle \geq 1 - \tilde{\eta}_1$ , while for all other  $j \neq l$ ,  $d^{-1}\langle Y, \tilde{X}_l \rangle < 1 - \tilde{\eta}_2$ . If no such codeword exists, it returns "#". Given that message i is sent, the decoder succeeds upon the following event:

- If  $i \in i([m])$ , then: 1)  $d^{-1}\langle Y, \tilde{X}_{i^{-1}(i)} \rangle \ge 1 \tilde{\eta}_1$ ; 2) For all  $l \ne i^{-1}(i)$ ,  $d^{-1}\langle Y, \tilde{X}_l \rangle < 1 \tilde{\eta}_2$ .
- If  $i \notin i([m])$ , then for all  $l \in [m]$ ,  $d^{-1}\langle Y, \tilde{X}_l \rangle < 1 \tilde{\eta}_2$ .

From now on, we assume that  $i \in i([m])$ ; when  $i \notin i([m])$ , the analysis follows in a similar manner. Set  $e_l = \tilde{X}_l - X_{i(l)}$ , which may be chosen adversarially (but is independent of the noise Z), such that  $||e_l|| \le \sqrt{\varepsilon d}$  (by definition of  $\tilde{X} \in \operatorname{Approx}_{\varepsilon}(\mathcal{X}_k)$ ). We have

$$\begin{split} d^{-1}\langle \boldsymbol{Y}, \tilde{\boldsymbol{X}}_l \rangle &= d^{-1}\langle \boldsymbol{Y}, \boldsymbol{X}_{\mathsf{i}(l)} \rangle + d^{-1}\langle \boldsymbol{Y}, \boldsymbol{e}_l \rangle \\ &= d^{-1}\langle \boldsymbol{Y}, \boldsymbol{X}_{\mathsf{i}(l)} \rangle + d^{-1}\langle \boldsymbol{X}_i, \boldsymbol{e}_l \rangle + d^{-1}\sigma\langle \boldsymbol{Z}, \boldsymbol{e}_l \rangle \,. \end{split}$$

Clearly,  $|d^{-1}\langle X_i, e_l\rangle| \leq \sqrt{\varepsilon}$  (Cauchy-Schwartz). Set  $M = \max_{l \in [m]} |d^{-1}\sigma\langle Z, e_l\rangle|$ . Observe:

$$d^{-1}\langle \mathbf{Y}, \mathbf{X}_{i} \rangle \geq 1 - \eta_{1} \implies d^{-1}\langle \mathbf{Y}, \tilde{\mathbf{X}}_{i^{-1}(i)} \rangle \geq 1 - \eta_{1} - \sqrt{\varepsilon} - M,$$
  

$$d^{-1}\langle \mathbf{Y}, \mathbf{X}_{j} \rangle \leq 1 - \eta_{2} \implies d^{-1}\langle \mathbf{Y}, \tilde{\mathbf{X}}_{i^{-1}(j)} \rangle \leq 1 - \eta_{2} + \sqrt{\varepsilon} + M.$$
(73)

We claim that with probability  $1-2k^{-5}$ , it holds that  $M \leq C\sqrt{\varepsilon}$  for some  $C=C(\beta)$ . Consequently, if  $\varepsilon$  is small enough, and the thresholds  $\tilde{\eta}_1 < \tilde{\eta}_2$  are chosen such that

$$1 - \eta_1 - \sqrt{\varepsilon} - C\sqrt{\varepsilon} > 1 - \tilde{\eta}_1, \quad 1 - \eta_2 + \sqrt{\varepsilon} + C\sqrt{\varepsilon} < 1 - \tilde{\eta}_2,$$

then under the high-probability event  $S_i(\eta_1,\eta_2)\cap\{M\leq C\sqrt{\varepsilon}\}$ , the adapted decoder necessarily returns the correct message  $\mathsf{i}^{-1}(i)$ . Since  $1-\eta_1>1-\eta_2$ , then this can clearly be made to hold whenever  $\varepsilon<\varepsilon_0$  for small enough constant  $\varepsilon_0$ . The proof of Proposition 25 therefore concludes by the Lemma below.

**Lemma 26** Assume the conditions of Proposition 25 (with R = 0) and the setup described above. There is  $C = C(\beta)$  such that  $M \le C\sqrt{\varepsilon}$  holds with probability  $\ge 1 - 2k^{-5}$ .

**Proof** Observe that  $d^{-1}\sigma(\mathbf{Z}, \mathbf{e}_l)$  is Gaussian with mean 0 and variance  $\leq \varepsilon \frac{\sigma^2}{d}$ . and recall that, by definition  $M = \max_{l \in [m]} |d^{-1}\sigma(\mathbf{Z}, \mathbf{e}_l)|$ .

By standard results on the maxima of Gaussian random variables, Lemmas 41 and 42, there is some universal c such that  $M_i \leq \sqrt{\varepsilon \frac{\sigma^2}{d}} \cdot \sqrt{c \log k}$  holds with probability  $\geq 1 - 2k^{-5}$ .

It remains to observe that

$$\frac{\sigma^2}{d}\log k = \sigma^2 \mathsf{R}_{d,k} = \sigma^2 \mathsf{C}(\beta \sigma^2) \le \frac{1}{2\beta}\,,$$

where we used  $C(s) = \frac{1}{2} \log(1 + 1/s) \le 1/(2s)$ .

## D.2.2. Positive Rate (R > 0)

We adapt the decoder from Section A.2. Recall the decision rule implemented by this decoder (assuming access to the true codebook  $\mathcal{X}_k$ ): the decoder returns an index i if: 1)  $d^{-1/2} \| \alpha \mathbf{Y} - \mathbf{X}_i \| \leq \sqrt{\tau_1}$ ; 2) For all  $j \neq i$ ,  $d^{-1/2} \| \alpha \mathbf{Y} - \mathbf{X}_i \| > \sqrt{\tau_2}$ ; if no such i exists, it returns #. Here  $0 < \sqrt{\tau_1} < \sqrt{\tau_2}$  are appropriately chosen thresholds.

We now adapt the aforementioned decoder to use  $\tilde{\mathcal{X}}$  instead of  $\mathcal{X}_k$ .

Observe that by the triangle inequality,

$$d^{-1/2}\|\alpha \boldsymbol{Y} - \boldsymbol{X}_{\mathsf{i}(l)}\| - \sqrt{\varepsilon} \le d^{-1/2}\|\alpha \boldsymbol{Y} - \tilde{\boldsymbol{X}}_l\| \le d^{-1/2}\|\alpha \boldsymbol{Y} - \boldsymbol{X}_{\mathsf{i}(l)}\| + \sqrt{\varepsilon}.$$

Consequently,

$$d^{-1/2}\|\alpha \mathbf{Y} - \mathbf{X}_{\mathsf{i}(l)}\| \le \sqrt{\tau_1} \implies d^{-1/2}\|\alpha \mathbf{Y} - \tilde{\mathbf{X}}_l\| \le \sqrt{\tau_1} + \sqrt{\varepsilon},$$
  

$$d^{-1/2}\|\alpha \mathbf{Y} - \mathbf{X}_{\mathsf{i}(l)}\| > \sqrt{\tau_2} \implies d^{-1/2}\|\alpha \mathbf{Y} - \tilde{\mathbf{X}}_l\| > \sqrt{\tau_2} - \sqrt{\varepsilon}.$$
(74)

The adapted decoder will operate as follows. Assume that  $\varepsilon < \varepsilon_0$  for some  $2\sqrt{\varepsilon_0} < \sqrt{\tau_2} - \sqrt{\tau_1}$ , and set

$$\sqrt{\tilde{\tau}_1} = \sqrt{\tau_1} + \sqrt{\varepsilon_0}, \quad \sqrt{\tilde{\tau}_2} = \sqrt{\tau_2} - \sqrt{\varepsilon_0}.$$

The adapted decoder implements the following rule. It returns an index  $l \in [m]$  whenever:  $\tilde{\boldsymbol{X}}_l$  is such that  $\|\alpha \boldsymbol{Y} - \tilde{\boldsymbol{X}}_l\| \leq \sqrt{\tilde{\tau}_1}$ ; 2) For all  $j \in [m] \setminus \{l\}$ ,  $\|\alpha \boldsymbol{Y} - \tilde{\boldsymbol{X}}_j\| > \sqrt{\tilde{\tau}_2}$ ; if no such l exists, it returns #.

Let us bound the error probability of the decoder. Suppose that a message i was sent, so that  $Y = X_i + \sigma Z$ , and assume without loss of generality that  $i \in i([m])$ ; the case  $i \notin i([m])$  follows similarly. Consider the event

$$S(\tau_1, \tau_2) = \{d^{-1/2} \|\alpha \mathbf{Y} - \mathbf{X}_i\| \le \sqrt{\tau_1}\} \cap \bigcap_{j \in [k] \setminus \{i\}} \{d^{-1/2} \|\alpha \mathbf{Y} - \mathbf{X}_j\| > \sqrt{\tau_2}\}.$$

In Section A.2 it was shown that, for appropriately chosen  $\tau_1, \tau_2, S(\tau_1, \tau_2)$  is a high-probability event. Notice that by (74), under the event  $S(\tau_1, \tau_2)$ , the adapted decoder indeed returns  $i^{-1}(i)$ ; thus, we have proven that the average error probability is vanishing.

### D.3. Proof of Lemma 12

Towards the proof of Lemma 12, we analyze the performance of Step II of the algorithm (Section 4.2) under a slightly more general setting, that does not use the strong symmetry properties that are available (w.h.p.) for a random spherical codebook  $\mathcal{X}_k$ , and that were crucial for constructing the test of Step I (Section 4.1). Instead, we only assume that one has access to some decoder under which the codebook attains small average error probability.

Let  $\mathcal{X}_k \subseteq (\sqrt{d}\mathcal{S}^{d-1})^k$  be a fixed codebook. Let  $\Psi: \mathbb{R}^d \to [k] \cup \{\#\}$  be a decoder. Denote by  $P_{i,j} = \Pr_{\mathbf{Z}}(\Psi(\mathbf{X}_i + \sigma \mathbf{Z}) = j)$  the probability (over the noise  $\mathbf{Z}$ ) that  $\Psi(\cdot)$  outputs symbol  $j \in [k] \cup \{\#\}$  given that the true label was i. To keep the presentation light, we start by introducing some notation.

We say that the decoder  $\Psi(\cdot)$  satisfies the  $(\rho, \varphi)$ - average error probability guarantee if there exists an index set  $\mathcal{I} \subseteq [k]$  of size  $|\mathcal{I}| \geq (1-\varphi)k$  such that  $\mathrm{Range}(\Psi) \subseteq \mathcal{I} \cup \{\#\}$  and the following holds. Denote by

$$\bar{P}_i = \begin{cases} 1 - P_{i,i} & \text{if } i \in \mathcal{I}, \\ 1 - P_{i,\#} & \text{if } i \notin \mathcal{I} \end{cases}$$

$$(75)$$

the *error probability* of the i-th message. Note that this is the same notion of error probability as in (71) from Section D.1 above. Then we have

$$\frac{1}{k} \sum_{i=1}^{k} \bar{P}_i \le \rho. \tag{76}$$

Note: under the conditions of Lemma 12, upon successful completion of Step I, and with high probability over  $\mathcal{X}_k \sim \mathrm{Unif}(\sqrt{d}\mathcal{S}^{d-1})^{\otimes k}$ , Proposition 25 implies that one may construct a decoder  $\Psi(\cdot) = \mathrm{Dec}(\cdot|\tilde{\mathcal{X}}_{\mathrm{I}})$  which satsifes the  $(\rho,\varphi)$  average error probability guarantee (up to a global relabeling, which we shall ignore henceforth), for some  $\rho = o(1)$ .

Let us describe once again in detail the procedure of Step II, stated in terms of the notation above. One has access to a decoder (constructed from  $\tilde{\mathcal{X}}_{\mathrm{I}}$ ), and uses it to label a batch of  $\bar{N}$  new samples,  $\mathbf{Y}_{1},\ldots,\mathbf{Y}_{\bar{N}}$ . Let  $S_{i}\subseteq[\bar{N}]$  be the subset of all measurements that have been assigned label i:

$$S_i = \{j : \Psi(\mathbf{Y}_j) = i\} .$$

Note that measurements assigned # are simply discarded. Next, we compute the cluster means:

$$\boldsymbol{A}_i = \frac{1}{|S_i|} \sum_{j \in S_i} \boldsymbol{Y}_j \,,$$

so that  $A_i = 0$  if  $S_i = \emptyset$ . The final centers returned by the procedure are simply the projections of  $A_i$  onto the ball  $\mathcal{B}(\mathbf{0}, \sqrt{d})$ :<sup>4</sup>

$$\hat{\boldsymbol{X}}_i = \mathcal{P}(\boldsymbol{A}_i)$$
 .

The next Lemma summarizes our guarantees for Step II as described above.

<sup>4.</sup> Note that we project onto the ball rather than the sphere  $\sqrt{d}S^{d-1}$  since projection onto convex sets is contracting in Euclidean norm. This is not so much the case for projection onto the sphere.

**Lemma 27** Suppose that  $\beta > 1$ , and let  $\varepsilon, \varphi \in (0,1)$  be constants. Let  $\mathcal{X}_k \subset (\sqrt{d}\mathcal{S}^{d-1})^k$  be some fixed center configuration, and suppose that one is given a decoder  $\Psi(\cdot)$  satisfying the  $(\rho, \varphi)$  average error probability guarantee, for some arbitrary  $\rho = o_{\beta,R}(1)$ .

Suppose that Step II is run with  $\bar{N} \geq \frac{k\sigma^2}{\varepsilon} + C\frac{k}{\varepsilon^{1/2}}\log(1/\varphi)$  for some sufficiently large univeral C>0. Then

 $\mathbb{E}\mathcal{L}_{avg}(\boldsymbol{\mathcal{X}}_k, \hat{\boldsymbol{\mathcal{X}}}_k) \leq \frac{\varepsilon}{1 - \varepsilon^{1/4}} + 8\varphi + o_{\beta,\mathsf{R}}(1)\,,$ 

where the expectation is only taken over the randomness in the measurements  $\mathbf{Y}_j = \mathbf{X}_{\ell_j} + \mathbf{Z}_j$ , namely  $\ell_j \sim \mathrm{Unif}([k])$  and  $\mathbf{Z}_j \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  (and the rate of decay in the o(1) term depends on  $\rho$ ).

The proof of Lemma 27 shall be given momentarily, in Section D.4 below. Before getting to it, let us show how it immediately implies Lemma 12:

**Proof** (Of Lemma 12). Note that the diameter of the ball  $\mathcal{B}(\mathbf{0}, \sqrt{d})$  is  $2\sqrt{d}$ , so necessarily  $\mathrm{dist}^2(\boldsymbol{X}_i, \hat{\boldsymbol{\mathcal{X}}}_k) \leq 4d$  and therefore  $\mathcal{L}_{\mathrm{avg}}(\boldsymbol{\mathcal{X}}_k, \hat{\boldsymbol{\mathcal{X}}}_k) \leq 4$ . Thus, for any event  $\mathcal{E}$ ,

$$\mathbb{E}\mathcal{L}_{\text{avg}}(\boldsymbol{\mathcal{X}}_k, \hat{\boldsymbol{\mathcal{X}}}_k) \leq \mathbb{E}\left[\mathcal{L}_{\text{avg}}(\boldsymbol{\mathcal{X}}_k, \hat{\boldsymbol{\mathcal{X}}}_k)\mathbb{1}\{\mathcal{E}\}\right] + 4\Pr(\mathcal{E}^c). \tag{77}$$

Let  $\tilde{\mathcal{X}}_I$  be the list returned by Step I of the algorithm. Let  $\varepsilon_0 = \varepsilon_0(\beta, \mathsf{R})$  be the threshold from Proposition 25, and consider the event

$$\mathcal{E}_I = \left\{ \tilde{\boldsymbol{\mathcal{X}}}_I \in \operatorname{Approx}_{\varepsilon_0}(\boldsymbol{\mathcal{X}}_k), \ |\mathcal{I}| \ge (1 - \varphi)k \right\},$$

where  $\mathcal{I}$  is defined in Lemma 11. By Lemma 11, provided that N is large enough,  $\Pr(\mathcal{E}_1^c) \leq \varphi + o_{\beta,R}(1)$ . Let  $\Psi(\cdot) = \mathsf{Dec}(\cdot|\tilde{\boldsymbol{\mathcal{X}}}_1)$  be the decoder promised by Proposition 25. For a suitably chosen (large enough)  $\rho = o_{\beta,R}(1)$ , denote the event

$$\mathcal{E} = \{\Psi(\cdot) \text{ satisfies the } (\rho, \varphi) \text{ average probability guarantee} \}$$
 .

By Proposition 25, along with Markov's inequality, we have  $\Pr(\mathcal{E}^c|\mathcal{E}_I) = o_{\beta,R}(1)$ . Thus,

$$\Pr(\mathcal{E}^c) \le \Pr(\mathcal{E}_I^c) + \Pr(\mathcal{E}^c | \mathcal{E}_I) \le \varphi + o_{\beta, R}(1)$$
.

Use the event  $\mathcal{E}$  in (77). By Lemma 27, for large enough  $\bar{N}$ ,  $\mathbb{E}\left[\mathcal{L}_{avg}(\boldsymbol{\mathcal{X}}_k, \hat{\boldsymbol{\mathcal{X}}}_k)\mathbb{1}\{\mathcal{E}\}\right] \leq \frac{\varepsilon}{1-\varepsilon^{1/4}} + 8\varphi + o_{\beta,R}(1)$ , and so Lemma 12 follows.

#### D.4. Proof of Lemma 27

Conceptually, the proof of Lemma 27 is quite straightforward. Denote by  $L_i \subseteq [\bar{N}]$  the measurements  $Y_j = X_{\ell_j} + Z_j$  whose true label is  $\ell_j = i$ . Imagine, for a moment, that we had access to a genie-aided decoder, that always assigns measurements to their true labels. In that case,  $S_i = L_i$ , and so  $A_i$  is just the sample mean of  $|L_i|$  i.i.d. Gaussian measurements  $\mathcal{N}(X_i, \sigma^2 I)$ . Since, on average,  $|L_i| = \bar{N}/k$ , the MSE is  $\mathbb{E}||X_i - A_i||^2 \approx \sigma^2 k/\bar{N}$ , which is  $\approx \varepsilon$  when  $\bar{N} \approx \sigma^2 k/\varepsilon$ . In practice, however, one does not have access to a clairvoyant decoder: we only assume an average error probability guarantee. Another difficulty is that we cannot guarantee that  $|L_i| \approx \bar{N}/k$  simultaneously for all i, unless  $\sigma^2 \gtrsim \log k$  (recall the coupon-collecting issue highlighted in the main

paper: we need  $\bar{N} \gtrsim k \log k$  to even *observe* a measurement of every label). Consequently, the analysis has to be carried out somewhat delicately.

We start with the trivial observation, that the average error probability guarantee implies that, in fact, *most* individual labels  $i \in [k]$  must have a small error probability:

**Lemma 28** Let  $(\mathcal{X}_k, \Psi(\cdot))$  satisfy the  $(\rho, \varphi)$  average error probability guarantee (76). There is a set of indices  $\mathcal{J}_1 \subseteq [k]$  of size  $|\mathcal{J}_1| \ge (1 - \rho^{1/2})k$  such that

$$\bar{P}_i < \rho^{1/2} \quad \text{for all } i \in \mathcal{J}_1 \,.$$
 (78)

**Proof** An immediate consequence of Markov's inequality.

Aside from lower bounding  $|L_i| \gtrsim \bar{N}/k$  in expectation, we shall also need to control the number of measurements that were *erroneously* assigned label i by  $\Psi(\cdot)$ . Following the notation of (75), let

$$\bar{Q}_i = \frac{1}{k} \sum_{l \in [k] \setminus \{i\}} P_{l,i} \tag{79}$$

be the probability that a random measurement  $Y = X_{\ell} + \sigma Z$ ,  $\ell \sim \text{Unif}([k])$  has true label  $\ell \neq i$ , but is erroneously assigned label i by  $\Psi(\cdot)$ .

**Lemma 29** Let  $(\mathcal{X}_k, \Psi(\cdot))$  satisfy the  $(\rho, \varphi)$  average error probability guarantee (76). There is a set of indices  $\mathcal{J}_2 \subseteq [k]$  of size  $|\mathcal{J}_2| \ge (1 - \rho^{1/2})k$  such that

$$\bar{Q}_i \le \rho^{1/2}/k \quad \text{for all } i \in \mathcal{J}_2 \,.$$
 (80)

**Proof** Observe that

$$\frac{1}{k} \sum_{i=1}^{k} \left( \sum_{l \in [k] \setminus \{i\}} P_{l,i} \right) = \frac{1}{k} \sum_{l=1}^{k} \left( \sum_{i \in [k] \setminus \{l\}} P_{l,i} \right) \le \frac{1}{k} \sum_{l=1}^{k} \bar{P}_{l} \le \rho.$$

Consequently, by Markov's inequality, there is  $\mathcal{J}_2 \subseteq [k]$  of size  $|\mathcal{J}_2| \geq (1 - \rho^{1/2})k$  such that  $\sum_{l \in [k] \setminus \{i\}} P_{l,i} \leq \rho^{1/2}$  for all  $i \in \mathcal{J}_2$ .

Recall:  $L_i \subseteq [\bar{N}]$  are the measurements whose true label is i;  $S_i \subseteq [\bar{N}]$  are the measurements assigned label i by  $\Psi(\cdot)$  (whether truthfully or erroneously). Define the event

$$\mathcal{E}_{i,1} = \left\{ |S_i \setminus L_i| \le \rho^{1/4} \frac{\bar{N}}{k} \quad \cap \quad |L_i| \ge (1 - \varepsilon^{1/4}) \frac{\bar{N}}{k} \right\}. \tag{81}$$

**Lemma 30** *Under the conditions of Lemma 27, for any*  $i \in \mathcal{J}_2$ *,* 

$$\Pr\left(\mathcal{E}_{i,1}^c\right) \le \varphi + o(1)$$
.

**Proof** Let us start by showing that  $|S_i \setminus L_i|$  is small with high probability. By definition,  $\Pr(\mathbf{Y} \in S_i \setminus L_i) = \bar{Q}_i$ . Since  $i \in \mathcal{J}_2$  this probability is  $\leq \rho^{1/2}/k$ . Thus, by Markov's inequality,

$$\Pr\left(|S_i \setminus L_i| > \rho^{1/4} \frac{\bar{N}}{k}\right) \le \frac{\mathbb{E}|S_i \setminus L_i|}{\rho^{1/4} \frac{\bar{N}}{k}} \le \frac{\rho^{1/2} \frac{\bar{N}}{k}}{\rho^{1/4} \frac{\bar{N}}{k}} = \rho^{1/4} = o(1).$$

Moving on, observe that  $|L_i| \sim \mathrm{Binomial}(1/k, \bar{N})$ . By Chernoff's inequality, Lemma 40,

$$\Pr\left(|L_i| \le (1 - \varepsilon^{1/4}) \frac{\bar{N}}{k}\right) \le e^{-c\varepsilon^{1/2} \frac{\bar{N}}{k}},$$

which is  $\leq \varphi$  for  $\bar{N} \geq C \frac{k}{\varepsilon^{1/2}} \log(1/\varphi)$ .

For  $i \in \mathcal{I}$ , define the event

$$\mathcal{E}_{i,2} = \left\{ |S_i \cap L_i| \ge (1 - \rho^{1/4})|L_i| \right\}. \tag{82}$$

**Lemma 31** *Under the conditions of Lemma 27, for any*  $i \in \mathcal{I} \cap \mathcal{J}_1$ *,* 

$$\Pr\left(\mathcal{E}_{i,2}^c\right) = o(1).$$

**Proof** By the definition of  $\mathcal{J}_1$ , conditioned on  $\ell_j = i$ ,  $\Pr(j \notin S_i | \ell_j = i) = \bar{P}_i \leq \rho^{1/2}$ . Thus, by Markov's inequality,

 $\Pr(|L_i \setminus S_i| \ge \rho^{1/4} |L_i|) \le \rho^{1/4} = o(1).$ 

We are ready to bound  $\mathbb{E}\mathcal{L}_{avg}(\mathcal{X}_k, \hat{\mathcal{X}}_k)$ . Observe that for any events  $\mathcal{E}_1, \dots, \mathcal{E}_k$ ,

$$\mathbb{E}\mathcal{L}_{avg}(\boldsymbol{\mathcal{X}}_{k}, \hat{\boldsymbol{\mathcal{X}}}_{k}) = \frac{1}{k} \sum_{i=1}^{k} d^{-1}\mathbb{E}\operatorname{dist}^{2}(\boldsymbol{X}_{i}, \hat{\boldsymbol{\mathcal{X}}}_{k})$$

$$\stackrel{(i)}{\leq} \frac{1}{k} \sum_{i=1}^{k} d^{-1}\mathbb{E}\left[\operatorname{dist}^{2}(\boldsymbol{X}_{i}, \hat{\boldsymbol{\mathcal{X}}}_{k})\mathbb{1}\left\{\mathcal{E}_{i}, i \in \mathcal{I} \cap \mathcal{J}_{1} \cap \mathcal{J}_{2}\right\}\right]$$

$$+ \frac{4}{k} \sum_{i=1}^{k} \Pr(\mathcal{E}_{i}^{c}) + 4\frac{|\mathcal{I}^{c}|}{k} + 4\frac{|(\mathcal{J}_{1} \cap \mathcal{J}_{2})^{c}|}{k}$$

$$\stackrel{(ii)}{\leq} \frac{1}{k} \sum_{i=1}^{k} d^{-1}\mathbb{E}\left[\operatorname{dist}^{2}(\boldsymbol{X}_{i}, \hat{\boldsymbol{\mathcal{X}}}_{k})\mathbb{1}\left\{\mathcal{E}_{i}, i \in \mathcal{I} \cap \mathcal{J}_{1} \cap \mathcal{J}_{2}\right\}\right] + \frac{4}{k} \sum_{i=1}^{k} \Pr(\mathcal{E}_{i}^{c}) + 4\varphi + o(1),$$

$$(83)$$

where: (i) follows from  $\operatorname{dist}(\boldsymbol{\mathcal{X}}_k, \hat{\boldsymbol{\mathcal{X}}}_k) \leq 4d$  (the diameter of the ball is  $2\sqrt{d}$ ); (ii) Follows from  $|\mathcal{I}| \leq \varphi k$  (by definition of the  $(\rho, \varphi)$  error probability guarantee) and  $|(\mathcal{J}_1 \cap \mathcal{J}_2)^c| = o(k)$  follows from Lemmas 28 and 29, with  $\rho = o(1)$ .

We take the event

$$\mathcal{E}_i = \mathcal{E}_{i,1} \cap \mathcal{E}_{i,2} \cap \mathcal{E}_{i,3} \,, \tag{84}$$

where  $\mathcal{E}_{i,1}$  is defined in (81),  $\mathcal{E}_{i,2}$  is defined in (82), and the definition of  $\mathcal{E}_{i,3}$  shall be deferred for later; its details would be somewhat obtuse at this point in the analysis.

To lighten the notation, introduce

$$\mathbb{E}_{\text{avg},\mathcal{E}}[F_i] = \frac{1}{kd} \sum_{i=1}^k \mathbb{E}[F_i \mathbb{1} \{\mathcal{E}_i, i \in \mathcal{I} \cap \mathcal{J}_1 \cap \mathcal{J}_2\}],$$
 (85)

where  $F_i$  is any sequence indexed by  $i \in [k]$ .

For  $i \in \mathcal{I}$ ,

$$\operatorname{dist}^{2}(X_{i}, \hat{X}_{k}) \leq ||X_{i} - \hat{X}_{i}||^{2} \leq ||X_{i} - A_{i}||^{2}$$

where we used  $\hat{X}_i = \mathcal{P}(A_i)$  and that projection onto convex sets is contracting with respect to Euclidean norm. Moreover, note that

$$\mathcal{D} = (D_1, \dots, D_k) \in \mathbb{R}^{d \times k} \quad \mapsto \quad (\mathbb{E}_{avg, \mathcal{E}}[\|D_i\|^2])^{1/2}$$

is a semi-norm, hence satisfies the triangle inequality.

We decompose

$$\begin{aligned} \boldsymbol{A}_i &= \frac{1}{|S_i|} \sum_{j \in S_i} \boldsymbol{Y}_j \\ &= \frac{1}{|S_i|} \sum_{j \in L_i} \boldsymbol{Y}_j + \frac{1}{|S_i|} \sum_{j \in S_i \setminus L_i} \boldsymbol{Y}_j - \frac{1}{|S_i|} \sum_{j \in L_i \setminus S_i} \boldsymbol{Y}_j \\ &= \frac{1}{|S_i|} \sum_{j \in L_i} \boldsymbol{Y}_i + \frac{1}{|S_i|} \sum_{j \in S_i \setminus L_i} \boldsymbol{X}_i - \frac{1}{|S_i|} \sum_{j \in L_i \setminus S_i} \boldsymbol{X}_{\ell_j} + \frac{\sigma}{|S_i|} \sum_{j \in S_i \setminus L_i} \boldsymbol{Z}_j - \frac{\sigma}{|S_i|} \sum_{j \in L_i \setminus S_i} \boldsymbol{Z}_j , \end{aligned}$$

and accordingly bound the first term of (83):

$$\left(\mathbb{E}_{\text{avg},\mathcal{E}}\text{dist}^{2}(\boldsymbol{X}_{i}, \hat{\boldsymbol{X}}_{k})\right)^{1/2} \leq \left(\mathbb{E}_{\text{avg},\mathcal{E}} \|\boldsymbol{X}_{i} - \boldsymbol{A}_{i}\|^{2}\right)^{1/2} \\
\leq \underbrace{\left(\mathbb{E}_{\text{avg},\mathcal{E}} \left\|\boldsymbol{X}_{i} - \frac{|L_{i}|}{|S_{i}|} \boldsymbol{X}_{i}\right\|^{2}\right)^{1/2}}_{I_{1}} + \underbrace{\left(\mathbb{E}_{\text{avg},\mathcal{E}} \left\|\frac{|L_{i}|}{|S_{i}|} \boldsymbol{X}_{i} - \frac{1}{|S_{i}|} \sum_{j \in L_{i}} \boldsymbol{X}_{j}\right\|^{2}\right)^{1/2}}_{I_{2}} \\
+ \underbrace{\left(\mathbb{E}_{\text{avg},\mathcal{E}} \left\|\frac{1}{|S_{i}|} \sum_{j \in S_{i} \setminus L_{i}} \boldsymbol{X}_{i} - \frac{1}{|S_{i}|} \sum_{j \in L_{i} \setminus S_{i}} \boldsymbol{X}_{\ell_{j}}\right\|^{2}\right)^{1/2}}_{I_{3}} \\
+ \underbrace{\left(\mathbb{E}_{\text{avg},\mathcal{E}} \left\|\frac{\sigma}{|S_{i}|} \sum_{j \in S_{i} \setminus L_{i}} \boldsymbol{Z}_{j}\right\|^{2}\right)^{1/2}}_{I_{3}} + \underbrace{\left(\mathbb{E}_{\text{avg},\mathcal{E}} \left\|\frac{\sigma}{|S_{i}|} \sum_{j \in L_{i} \setminus S_{i}} \boldsymbol{Z}_{j}\right\|^{2}\right)^{1/2}}_{I_{5}}. \tag{86}$$

We proceed to bound the terms above. Starting with  $I_1$ , observe that

$$I_1^2 = \mathbb{E}_{ ext{avg},\mathcal{E}} \left\| oldsymbol{X}_i - rac{|L_i|}{|S_i|} oldsymbol{X}_i 
ight\|^2 = d \cdot \mathbb{E}_{ ext{avg},\mathcal{E}} \left| 1 - rac{|L_i|}{|S_i|} 
ight|,$$

and therefore  $I_1 = o(1)$ , since under  $\mathcal{E}_i$  we must have  $(1 - o(1))|L_i| \le |S_i| \le (1 + o(1))|L_i|$ . Moving on to  $I_2$ ,

$$I_2^2 = \mathbb{E}_{\mathbf{avg},\mathcal{E}} \left\| rac{|L_i|}{|S_i|} oldsymbol{X}_i - rac{1}{|S_i|} \sum_{j \in L_i} oldsymbol{Y}_j 
ight\|^2 = \mathbb{E}_{\mathbf{avg},\mathcal{E}} \left[ rac{\sigma^2 |L_i|}{|S_i|^2} \left\| rac{1}{\sqrt{|L_i|}} \sum_{j \in L_i} oldsymbol{Z}_j 
ight\|^2 
ight] \,.$$

Noting that  $\frac{1}{\sqrt{|L_i|}}\sum_{j\in L_i}\mathbf{Z}_j\sim\mathcal{N}(\mathbf{0},\mathbf{I})$ , and that under  $\mathcal{E}_i, \frac{|L_i|}{|S_i|^2}\leq (1+o(1))\frac{1}{1-\varepsilon^{1/4}}\frac{k}{N}$  we get that  $I_2^2\leq \frac{\varepsilon}{1-\varepsilon^{1/4}}+o(1)$  since  $\bar{N}\geq \sigma^2k/\varepsilon$ . As for  $I_3$ ,

$$I_3^2 = \mathbb{E}_{\mathbf{avg},\mathcal{E}} \left\| \frac{1}{|S_i|} \sum_{j \in S_i \setminus L_i} \boldsymbol{X}_i - \frac{1}{|S_i|} \sum_{j \in L_i \setminus S_i} \boldsymbol{X}_{\ell_j} \right\|^2 \leq d \cdot \mathbb{E}_{\mathbf{avg},\mathcal{E}} \left( \frac{|S_i \setminus L_i| + |L_i \setminus S_i|}{|S_i|} \right)^2 = o(1) .$$

Bounding the terms  $I_4$ ,  $I_5$  is somewhat more involved, and requires the introduction of a new event,  $\mathcal{E}_{i,3}$ , whose definition has been deferred up to this point.

**Lemma 32** Define the event  $\mathcal{E}_{i,3}$  in (88). Then  $\Pr(\mathcal{E}_{i,3}) = o(1)$  and moreover

$$I_4, I_5 = o(1)$$
.

To keep the narrative flow, we defer the proof of Lemma 32 to Section D.4.1 below.

We are ready to tie all loose ends, and finish the proof of Lemma 27. By (86) and the upper bounds we have shown for  $I_1, \ldots, I_5$ , we get  $\mathbb{E}_{\text{avg},\mathcal{E}} \text{dist}^2(\boldsymbol{X}_i, \hat{\boldsymbol{X}}_k) \leq \frac{\varepsilon}{1-\varepsilon^{1/4}} + o(1)$ . Combining Lemmas 30, 31 and 32, we get  $\Pr(\mathcal{E}_i) \leq \varphi + o(1)$ . Plugging these estimates into (83), we finally get the claimed bound of Lemma 27.

#### D.4.1. Proof of Lemma 32

The terms  $I_4$ ,  $I_5$  correspond to sums of independent "noise" vectors, whose mean is zero. Therefore, one expects different  $\mathbf{Z}_j$ -s to cancel out one another on average, so that, for example (considering  $I_4$ ),

$$\mathbb{E} \left\| \frac{\sigma}{|S_i|} \sum_{j \in S_i \setminus L_i} \mathbf{Z}_j \right\|^2 \approx \frac{\sigma^2}{|S_i|^2} \cdot |S_i \setminus L_i| d$$

rather than  $\frac{\sigma^2}{|S_i|^2} \cdot |S_i \setminus L_i|^2 d$  which is what we would have gotten by naive application of the triangle inequality. A subtle point is that the set  $S_i$  actually depends on the noise vectors  $\mathbf{Z}_j$ , so one needs to apply some care when taking the expectation above. We propose to overcome this difficulty through a rather crude bound.

<sup>5.</sup> A priori, we cannot discount the possibility that conditioned on  $j \in S_i$ , the noise  $\mathbb{Z}_j$  biases towards some particular direction.

For a subset  $B\subseteq [n]$  (which itself may be random, but independent of  $\{{\pmb Z}_j\}_{j\in [\bar N]}$ , let

$$D(B,m) = \max_{S \subseteq B, |S| \le m} \left\| \sum_{j \in S} Z_j \right\|^2.$$

Since under  $\mathcal{E}_{i,1} \cap \mathcal{E}_{i,2} \subset \mathcal{E}_i$  we have

$$|S_i \setminus L_i| \le \rho^{1/4} \frac{\bar{N}}{k}, \quad |L_i \setminus S_i| \le \rho^{1/4} |L_i|, \quad |L_i| \ge (1 - \varepsilon^{1/4}) \frac{\bar{N}}{k},$$

recalling the definition of  $I_4$ ,  $I_5$ , (86), clearly,

$$I_{4}^{2} = \sigma^{2} \cdot \mathbb{E}_{avg,\mathcal{E}} \left[ \frac{1}{|S_{i}|^{2}} \left\| \sum_{j \in S_{i} \setminus L_{i}} \mathbf{Z}_{j} \right\|^{2} \right]$$

$$\lesssim \frac{\sigma^{2}}{(\bar{N}/k)^{2}} \cdot \mathbb{E}_{avg,\mathcal{E}} \left[ \mathbf{D}([\bar{N}], \rho^{1/4}(\bar{N}/k)) \right],$$

$$I_{5}^{2} = \sigma^{2} \cdot \mathbb{E}_{avg,\mathcal{E}} \left[ \frac{1}{|S_{i}|^{2}} \left\| \sum_{j \in L_{i} \setminus S_{i}} \mathbf{Z}_{j} \right\|^{2} \right]$$

$$\lesssim \sigma^{2} \cdot \mathbb{E}_{avg,\mathcal{E}} \left[ \frac{1}{|L_{i}|^{2}} \cdot \mathbf{D}(L_{i}, \rho^{1/4} |L_{i}|) \right].$$
(87)

We are ready to define the event  $\mathcal{E}_{i,3}$ , which has been deferred up to this point.

The event  $\mathcal{E}_{i,3}$ . For C a sufficiently large universal constant, define

$$\mathcal{E}_{i,3} = \left\{ \mathbf{D}([n], \rho^{1/4}(\bar{N}/k)) \le C \left( \rho^{1/4}(\bar{N}/k) \right)^2 \log \frac{\bar{N}e}{\rho^{1/4}(\bar{N}/k)} + C\rho^{1/4}(\bar{N}/k) \log d \right\},$$

$$\cap \left\{ \mathbf{D}(L_i, \rho^{1/4}|L_i|) \le C \left( \rho^{1/4}|L_i| \right)^2 \log \frac{|L_i|e}{\rho^{1/4}|L_i|} + C\rho^{1/4}|L_i| \log d \right\}.$$
(88)

By Lemma 33, given below, C may indeed be chosen so that  $\Pr(\mathcal{E}_{i,3}^c) = o(1)$ .

**Bounding**  $I_4$ . Using (87) and (88),

$$I_4^2 \lesssim \frac{\sigma^2}{(\bar{N}/k)^2} \cdot d^{-1} \cdot \left\{ o\left((\bar{N}/k)^2 \log k\right) + o\left((\bar{N}/k) \log d\right) \right\} .$$

The first term is o(1) because  $\sigma^2 d^{-1} \log k = O(1)$ . The second term is o(1) because  $\frac{\sigma^2}{(\bar{N}/k)} d^{-1} \log d = O(d^{-1} \log d) = o(1)$ , since  $\bar{N} \gtrsim k \sigma^2$ . Thus  $I_4 = o(1)$ .

**Bounding**  $I_5$ . Using (87) and (88),

$$I_5^2 \lesssim \sigma^2 \mathbb{E}_{\text{avg},\mathcal{E}} \left[ o(1) + o\left(\frac{1}{|L_i|} \log d\right) \right]$$
  
 
$$\lesssim o\left(\sigma^2 d^{-1}\right) + o\left(\sigma^2 d^{-1} \frac{1}{(\bar{N}/k)} \log d\right).$$

The first term is o(1) since, by assumption,  $\sigma^2 = o(d)$ . Since  $\bar{N} \gtrsim k\sigma^2$ , the second terms is  $o\left(\frac{\log d}{d}\right) = o(1)$ . Thus,  $I_5 = o(1)$ .

This conclude the proof of Lemma 32.

## A Technical Lemma.

**Lemma 33** Let  $Z_1, ..., Z_n \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  be independent. For a set  $S \subseteq [n]$  let  $W_S = \sum_{i \in S} Z_i$ . There is a univeral C > 0 such that for  $t \ge 1$ ,

$$\Pr\left(\max_{S\subseteq[n],|S|\leq t}\|\boldsymbol{W}_S\|\geq Ct\sqrt{\log\frac{ne}{t}}+C\sqrt{t\log d}\right)\leq d^{-5}.$$

**Proof** This is a straightforward application of the well-known Lemmas 41 and Lemma 42, along with a standard "trick".

Let  $\mathcal{T}$  be a 1/2-net of  $\mathcal{S}^{d-1}$ , of size  $\leq 5^d$  (e.g. (Wainwright, 2019, Example 5.8)). It may be readily verified that for any vector  $\boldsymbol{x} \in \mathbb{R}^d$ ,  $\|\boldsymbol{x}\| = \max_{\boldsymbol{e} \in \mathcal{S}^{d-1}} \langle \boldsymbol{x}, \boldsymbol{e} \rangle \leq 2 \max_{\boldsymbol{e} \in \mathcal{T}} \langle \boldsymbol{x}, \boldsymbol{e} \rangle$ . Set  $N_t = |\mathcal{T}| \sum_{s=1}^t \binom{n}{s} \lesssim 5^d (ne/t)^t$ . Note also that for any  $\boldsymbol{e} \in \mathcal{S}^{d-1}$ ,  $\mathbb{E}\left[\langle \boldsymbol{W}_S, \boldsymbol{e} \rangle^2\right] \leq |S| \leq t$ .

By the expectation bound Lemma 41, and the Borell-TIS inequality, Lemma 42, for  $x \ge 0$ ,

$$\Pr\left(\max_{S\subseteq[n],|S|\leq t}\|\boldsymbol{W}_S\|\geq 2\sqrt{t}(\sqrt{2\log N_t}+x)\right)\leq \Pr\left(\max_{S\subseteq[n],|S|\leq t,\boldsymbol{e}\in\mathcal{T}}\langle\boldsymbol{W}_S,\boldsymbol{e}\rangle\geq \sqrt{t}(\sqrt{2\log N_t}+x)\right)$$
$$<\boldsymbol{e}^{-x^2/2}.$$

Set  $x = \sqrt{10 \log(d)}$  to get the claimed bound.

# **Appendix E. Auxiliary Technical Results**

### **E.1.** Concentration Inequalities

The Gaussian Lipschitz concentration inequality (Wainwright, 2019, Theorem 2.25):

**Lemma 34 (Gaussian Lipschitz concentration inequality)** *Let*  $f : \mathbb{R}^d \to \mathbb{R}$  *be L-Lipschitz, and*  $Z \sim \mathcal{N}(\mathbf{0}, I)$ . *For all*  $t \geq 0$ ,

$$\Pr(f(\mathbf{Z}) \ge \mathbb{E}f(\mathbf{Z}) + t) \le e^{-\frac{t^2}{2L^2}}.$$

Standard bound on the measure of a spherical cap (Wainwright, 2019, Eq. (3.33)):

**Lemma 35** Let  $\mathbf{Z} \sim \mathrm{Unif}(\mathcal{S}^{d-1})$ . For all  $\mathbf{u} \in \mathcal{S}^{d-1}$  and  $t \in (0,1)$ ,

$$\Pr\left(\langle \boldsymbol{u}, \boldsymbol{Z} \rangle \ge t\right) \le e^{-dt^2/2}$$
.

We state Bernstein's inequality for independent bounded random variables (Vershynin, 2018, Theorem 2.8.4)

**Lemma 36** Let  $X_1, \ldots, X_n$  be independent, centered, with  $|X_i| \leq K$ . Set  $S_n = \sum_{i=1}^n X_i$ . Then for all  $t \geq 0$ ,

$$\Pr(|S_n| \ge t) \le 2 \exp\left(-\frac{t^2/2}{\sum_{i=1}^n \text{Var}(X_i) + Kt/3}\right).$$

The following is Bernstein's inequality for sums of independent sub-Exponential random variables (Vershynin, 2018, Theorem 2.8.1):

**Lemma 37 (Bernstein's inequality, sub-Exponential RVs)** Let  $X_1, ..., X_n$  be independent and sub-exponential. Set  $S_n = \sum_{i=1}^n X_i$ . Then for all  $t \ge 0$ ,

$$\Pr(|S_n - \mathbb{E}[S_n]| \ge t) \le 2 \exp\left[-c \min\left(\frac{t^2}{\sum_{i=1}^n \|X_i\|_{\psi_1}^2}, \frac{t}{\max_{1 \le i \le n} \|X_i\|_{\psi_1}}\right)\right],$$

where c > 0 is an absolute constant.

We state the following version of the Chernoff bound for Bernoulli random variables. For a citable reference, see for example (Boucheron et al., 2013, Section 2.2):

**Lemma 38 (Chernoff's inequality)** Let  $X_1, \ldots, X_n$  be i.i.d. Bernoulli random variables, with  $\mathbb{E}[X_i] = q$ . Let  $t \leq q \leq p$ . Then

$$\Pr\left(\sum_{i=1}^{n} X_i \ge p\right) \le e^{nD_{\mathrm{KL}}(p;q)}, \quad \Pr\left(\sum_{i=1}^{n} X_i \le t\right) \le e^{nD_{\mathrm{KL}}(t;q)},$$

where

$$D_{\text{KL}}(p;q) = p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q}$$

is the Kullback-Leibler divergence divergence between Ber(p) and Ber(q).

We shall use Lemma 38 with the following easy estimate:

**Lemma 39** Let 0 < q < p < 1/2. Then

$$D_{\mathrm{KL}}(p;q) \ge p \log \frac{p}{q} - 2p$$
.

**Proof** One may readily verify that  $\log(1-p) \ge -2p$  holds for  $0 \le p \le 1/2$ . Thus,

$$D_{\mathrm{KL}}(p;q) = p\log\frac{p}{q} + (1-p)\log\frac{1-p}{1-q} \geq p\log\frac{p}{q} + (1-p)\log\left(1-p\right) \stackrel{(\star)}{\geq} p\log\frac{p}{q} + \log(1-p) \geq p\log\frac{p}{q} - 2p\,,$$

where  $(\star)$  holds since  $\log(1-p)$  is negative.

The following is a version of Chernoff's inequality, specialized for small deviations, and taken from (Vershynin, 2018, Exercise 2.3.5)

**Lemma 40 (Chernoff's inequality; small deviations)** *In the setting of Lemma 38, for*  $\delta \in (0,1)$ *,* 

$$\Pr\left(\sum_{i=1}^{n} X_i \ge (1+\delta)qn\right) \le e^{-c\delta^2qn}, \quad \Pr\left(\sum_{i=1}^{n} X_i \le (1-\delta)qn\right) \le e^{-c\delta^2qn},$$

where c > 0 is universal.

#### E.2. Maxima of Gaussian Random Variables

We state two elementary results about the maximum of n Gaussian random variables.

**Lemma 41** Let  $Z = (Z_1, ..., Z_n)$  be a Gaussian random vector, such that  $\mathbb{E}[Z] = \mathbf{0}$  and  $\mathbb{E}[Z_i^2] \leq \sigma^2$  for all i. Then

$$\mathbb{E}[\max_{1 \le i \le n} Z_i] \le \sqrt{2\sigma^2 \log n}.$$

(When  $Z_1, \ldots, Z_n$  are uncorrelated, this is in fact tight to leading order. But we shall not use this stronger fact.)

**Proof** This is classical. For completeness, we give a one-line proof. For all  $\beta > 0$ ,

$$\mathbb{E}[\max_{1 \le i \le n} Z_i] \le \frac{1}{\beta} \mathbb{E} \log \sum_{i=1}^n e^{\beta Z_i} \stackrel{(\star)}{\le} \frac{1}{\beta} \log \mathbb{E} \sum_{i=1}^n e^{\beta Z_i} \le \frac{1}{\beta} (\log n + \frac{1}{2} \sigma^2 \beta^2),$$

where  $(\star)$  follows from Jensen's inequality. Now set  $\beta = \sqrt{2\sigma^2 \log n}$ .

The following is a special (easy) case of the Borell-TIS inequality, see e.g. (Adler and Taylor, 2009, Theorem 2.1.1). Alternatively, this follows immediately from the Gaussian Lipschitz concentration inequality, Lemma 34:

**Lemma 42 (Borell-TIS)** Let  $\mathbf{Z} = (Z_1, \dots, Z_n)$  be a Gaussian random vector with  $\mathbb{E}[\mathbf{Z}] = \mathbf{0}$ . Set  $\sigma^2 = \max_{1 \le i \le n} \mathbb{E}Z_i^2$ . Then for  $t \ge 0$ ,

$$\Pr(\max_{1 \le i \le n} Z_i \ge \mathbb{E}[\max_{1 \le i \le n} Z_i] + t) \le e^{-\frac{t^2}{2\sigma^2}}, \quad \Pr(\max_{1 \le i \le n} Z_i \le \mathbb{E}[\max_{1 \le i \le n} Z_i] - t) \le e^{-\frac{t^2}{2\sigma^2}}.$$

## **E.3. Results From Information Theory**

The following is Fano's inequality, see e.g. (Polyanskiy and Wu, 2014, Theorem 5.2).

**Lemma 43 (Fano's inequality)** Let  $(\ell, Z, \hat{\ell})$  be random variables such that  $\ell, \hat{\ell} \in [k]$ , and the Markov chain

$$\ell \longrightarrow Z \longrightarrow \hat{\ell}$$

holds. Denote  $p_e = (\Pr(\ell \neq \hat{\ell}))$ . Then

$$H(\ell|Z) \le h_b(p_e) + p_e \log(k-1) ,$$

where  $h_b(p) = p \log \frac{1}{p} + (1-p) \log \frac{1}{1-p}$  is the binary entropy function.

Lastly is the celebrated I-MMSE relation of (Guo et al., 2005, Theorem 2):

**Lemma 44 (I-MMSE)** Let  $X \in \mathbb{R}^d$  be any random vector with finite second moments,  $\mathbb{E}||X||^2 < \infty$ . Let  $Z \sim \mathcal{N}(\mathbf{0}, I)$  be independent of X, and denote  $Y(s) = \sqrt{s}X + Z$ . Then

$$\frac{d}{ds}I(\boldsymbol{X};\boldsymbol{Y}(s)) = \frac{1}{2}\mathbb{E}\left[\|\boldsymbol{X} - \mathbb{E}(\boldsymbol{X}|\boldsymbol{Y}(s))\|^2\right] =: \text{mmse}(s).$$