On the Sample Complexity of Predictive Sparse Coding

Nishant A. Mehta* and Alexander G. Gray

College of Computing Georgia Institute of Technology Atlanta, GA 30332, USA

Abstract

The goal of predictive sparse coding is to learn a representation of examples as sparse linear combinations of elements from a dictionary, such that a learned hypothesis linear in the new representation performs well on a predictive task. Predictive sparse coding algorithms recently have demonstrated impressive performance on a variety of supervised tasks, but their generalization properties have not been studied. We establish the first generalization error bounds for predictive sparse coding, covering two settings: 1) the overcomplete setting, where the number of features k exceeds the original dimensionality d; and 2) the high or infinite-dimensional setting, where only dimension-free bounds are useful. Both learning bounds intimately depend on stability properties of the learned sparse encoder, as measured on the training sample. Consequently, we first present a fundamental stability result for the LASSO, a result characterizing the stability of the sparse codes with respect to perturbations to the dictionary. In the overcomplete setting, we present an estimation error bound that decays as $\tilde{O}(\sqrt{dk/m})$ with respect to d and k. In the high or infinite-dimensional setting, we show a dimension-free bound that is $\tilde{O}(\sqrt{k^2s/m})$ with respect to k and k, where k is an upper bound on the number of non-zeros in the sparse code for any training data point.

Keywords: Statistical learning theory, Luckiness, Data-dependent complexity, Dictionary learning, Sparse coding, LASSO

1 Introduction

Learning architectures such as the support vector machine and other linear predictors enjoy strong theoretical properties (Steinwart and Christmann, 2008; Kakade et al., 2009), but a learning-theoretic understanding of many more complex learning architectures is lacking. Predictive methods based on sparse coding recently have emerged which simultaneously learn a data representation via a nonlinear encoding scheme and an estimator linear in that representation (Bradley and Bagnell, 2009; Mairal et al., 2012, 2009). A sparse coding representation $z \in \mathbb{R}^k$ of a data point $x \in \mathbb{R}^d$ is learned by representing x as a sparse linear combination of k atoms $D_j \in \mathbb{R}^d$ of a dictionary $D = (D_1, \ldots, D_k) \in \mathbb{R}^{d \times k}$. In the coding $x \approx \sum_{j=1}^k z_j D_j$, all but a few z_j are zero. Predictive sparse coding methods such as Mairal et al. (2012)'s task-driven dictionary learning

Predictive sparse coding methods such as Mairal et al. (2012)'s task-driven dictionary learning recently have achieved state-of-the-art results on many tasks, including the MNIST digits task. Whereas standard sparse coding minimizes an unsupervised, reconstructive ℓ_2 loss, predictive sparse

^{*}To whom correspondence should be addressed. Email: niche@cc.gatech.edu

coding seeks to minimize a supervised loss by optimizing a dictionary and a linear predictor that operates on encodings to that dictionary. There is much empirical evidence that sparse coding can provide good abstraction by finding higher-level representations which are useful in predictive tasks (Yu et al., 2009). Intuitively, the power of *prediction-driven* dictionaries is that they pack more atoms in parts of the representational space where the prediction task is more difficult. However, despite the empirical successes of predictive sparse coding methods, it is unknown how well they generalize in a theoretical sense.

In this work, we develop what to our knowledge are the first generalization error bounds for predictive sparse coding algorithms; in particular, we focus on ℓ_1 -regularized sparse coding. Maurer and Pontil (2010) and Vainsencher et al. (2011) previously established generalization bounds for the classical, reconstructive sparse coding setting. Extending their analysis to the predictive setting introduces certain difficulties related to the richness of the class of sparse encoders. Whereas in the reconstructive setting, this complexity can be controlled directly by exploiting the stability of the reconstruction error to dictionary perturbations, in the predictive setting it appears that the complexity hinges upon the stability of the sparse codes themselves to dictionary perturbations. This latter notion of stability is much harder to prove; moreover, it can be realized only with additional assumptions which depend on the dictionary, the data, and their interaction (see Theorem 1). Furthermore, when the assumptions hold for the learned dictionary and data, we also need to guarantee that the assumptions hold on a newly drawn sample.

Contributions We provide learning bounds for two core scenarios in predictive sparse coding: the *overcomplete setting* where the dictionary size, or number of learned features, k exceeds the ambient dimension d; and the infinite-dimensional setting where only dimension-free bounds are acceptable. Both bounds hold provided the size m of the training sample is large enough, where the critical size for the bounds to kick in depends on a certain notion of stability of the learned representation. The core contributions of this work are:

- 1. Under mild conditions, a stability bound for the LASSO (Tibshirani, 1996) under dictionary perturbations. (Theorem 1)
- 2. In the overcomplete setting, a learning bound that is essentially of order $\sqrt{\frac{dk}{m}} + \frac{\sqrt{s}}{\lambda \mu_s(D)}$, where each sparse code has at most s non-zero coordinates. The term $\frac{1}{\mu_s(D)}$ is the inverse s-incoherence (see Definition 1) and is roughly the worst condition number among all linear systems induced by taking s columns of s. (Theorem 3)
- 3. In the infinite-dimensional setting, a learning bound that is *independent* of the dimension of the data; this bound is essentially of order $\frac{1}{\mu_{2s}(D)}\sqrt{\frac{k^2s}{m}}$. (Theorem 4)

The stability of the sparse codes are absolutely crucial to this work. Proving that the notion of stability of contribution 1 holds is quite difficult because the LASSO is not strongly convex in general. Consequently, much of the technical difficulty of this work is owed to finding conditions under which the LASSO is stable under dictionary perturbations and proving that when these conditions hold with respect to the learned hypothesis and the training sample, they also hold with respect to a future sample.

For convenience, we have collected all of the various notation of this paper in a glossary in Appendix G.

The predictive sparse coding problem

Let P be a probability measure over $B_{\mathbb{R}^d} \times \mathcal{Y}$, the product of an input space $B_{\mathbb{R}^d}$ (the unit ball of \mathbb{R}^d) and a space \mathcal{Y} of univariate labels; examples of \mathcal{Y} include a bounded subset of \mathbb{R} for regression and $\{-1,1\}$ for classification. Let $\mathbf{z}=(z_1,\ldots,z_m)$ be a sample of m points drawn iid from P, where each labeled point z_i equals (x_i, y_i) for $x_i \in B_{\mathbb{R}^d}$ and $y_i \in \mathcal{Y}$. In the reconstructive setting, labels are not of interest and we can just as well consider an unlabeled sample \mathbf{x} of m points drawn iid from the marginal probability measure Π on $B_{\mathbb{R}^d}$.

The sparse coding problem is to represent each point x_i as a sparse linear combination of k basis vectors, or atoms D_1, \ldots, D_k . The atoms form the columns of a dictionary D living in a space of dictionaries $\mathcal{D} := (B_{\mathbb{R}^d})^k$, for $D_i = (D_i^1, \dots, D_i^d)^T$ in the unit ℓ_2 ball. An encoder φ_D can be used to frame ℓ_1 sparse coding:

$$\varphi_D(x) := \underset{z}{\arg\min} \|x - Dz\|_2^2 + \lambda \|z\|_1;$$
 (1)

hence, encoding x as $\varphi_D(x)$ amounts to solving a LASSO problem. The reconstructive ℓ_1 sparse coding objective is then

$$\min_{D \in \mathcal{D}} \mathsf{E}_{x \sim \Pi} \|x - D\varphi_D(x)\|_2^2 + \lambda \|\varphi_D(x)\|_1,$$

Generalization bounds for the empirical risk minimization (ERM) variant of this objective have been established. In the infinite-dimensional setting, Maurer and Pontil (2010) showed that with probability $1 - \delta$ over the training sample **x**:

$$\sup_{D \in \mathcal{D}} Pf_D - P_{\mathbf{x}} f_D \leq \frac{k}{\sqrt{m}} \left(\frac{14}{\lambda} + \frac{1}{2} \sqrt{\log\left(16m/\lambda^2\right)} \right) + \sqrt{\frac{\log(1/\delta)}{2m}}$$
 (2)

where $f_D(x) := \min_{z \in \mathbb{R}^k} \|x - Dz\|_2^2 + \lambda \|z\|_1$. This bound is independent of the dimension d and hence useful when $d \gg k$, as in general Hilbert spaces. They also showed a similar bound in the overcomplete setting where the k is replaced by \sqrt{dk} . Vainsencher et al. (2011) handled the overcomplete setting, producing a bound that is $O(\sqrt{dk/m})$ as well as fast rates of O(dk/m), with only logarithmic dependence on $\frac{1}{\lambda}$.

Predictive sparse coding, introduced by Mairal et al. (2012), minimizes a supervised loss with respect to a representation and an estimator linear in the representation. Let \mathcal{W} be a space of linear hypotheses with $\mathcal{W} := rB_{\mathbb{R}^k}$, the ball in \mathbb{R}^k scaled to radius r. A predictive sparse coding hypothesis function f is identified by $f = (D, w) \in \mathcal{D} \times \mathcal{W}$ and defined as $f(x) = \langle w, \varphi_D(x) \rangle$. The function class \mathcal{F} is the set of such hypotheses. The loss will be measured via $l: \mathcal{Y} \times \mathbb{R} \to [0, b]$, b>0, a bounded loss function that is L-Lipschitz in its second argument.

The predictive sparse coding objective is²

$$\min_{D \in \mathcal{D}, w \in \mathcal{W}} \mathsf{E}_{(x,y) \sim P} l(y, \langle w, \varphi_D(x) \rangle) + \frac{1}{r} \|w\|_2^2; \tag{3}$$

¹To see this, take Theorem 1.2 of Maurer and Pontil (2010) with $Y=\{y\in\mathbb{R}^k:\|y\|_1<\frac{1}{\lambda}\}$ and

 $[\]mathcal{T} = \{T : \mathbb{R}^k \to \mathbb{R}^d : ||Te_j|| \le 1, j \in [k]\}$, so that $||\mathcal{T}||_Y \le \frac{1}{\lambda}$.

While the focus of this work is (3), formally the predictive sparse coding framework admits swapping out the squared ℓ_2 norm regularizer on w for any other regularizer.

In this work, we analyze the ERM variant of (3):

$$\min_{D \in \mathcal{D}, w \in \mathcal{W}} \frac{1}{m} \sum_{i=1}^{m} l(y_i, \langle w, \varphi_D(x_i) \rangle) + \frac{1}{r} \|w\|_2^2.$$

$$\tag{4}$$

Because this objective is not convex and global optimizers are not known, a priori we cannot say whether the optimal hypothesis or a nearly optimal hypothesis will be returned by any learning algorithm. However, we can and will bet on certain sparsity-related stability properties holding with respect to the learned hypothesis and the training sample. Consequently, all the presented learning bounds will hold uniformly not over the set of all hypotheses but rather potentially much smaller random subclasses of hypotheses. Additionally, the presented bounds will be algorithm-independent³, although certainly algorithm design can influence the observed stability of the learned hypothesis and hence the best learning bound that applies a posteriori.

Encoder stability Defining the encoder (1) via the ℓ_1 sparsity-inducing regularizer (sparsifier) is just one way of designing an encoder. The choice of sparsifier seems to be pivotal both from an empirical perspective and a theoretical one. Bradley and Bagnell (2009) used a differentiable approximate sparsifier based on the Kullback-Leibler divergence (true sparsity may not result). The ℓ_1 sparsifier $\|\cdot\|_1$ is the most popular and notably is the tightest convex lower bound for the ℓ_0 "norm": $\|x\|_0 := |\{i: x_i \neq 0\}|$ (Fazel, 2002). Regrettably, from a stability perspective the ℓ_1 sparsifier is not well-behaved in general. Indeed, due to the lack of strict convexity, each x need not have a unique image under φ_D . It also is unclear how to analyze the class of mappings φ_D , parameterized by D, if the map changes drastically under small perturbations to D. Hence, we will begin by establishing sufficient conditions under which φ_D is stable under perturbations to D.

2 Conditions and main results

In this section, we develop several quantities that are central to the statement of the main results. Throughout this paper, let $[n] := \{1, \ldots, n\}$ for $n \in \mathbb{N}$. Also, for $t \in \mathbb{R}^k$, define $\mathrm{supp}(t) := \{i \in [k] : t_i \neq 0\}$.

Definition 1 (s-incoherence) For $s \in [k]$ and $D \in \mathcal{D}$, the s-incoherence $\mu_s(D)$ is defined as the square of the minimum singular value among s-atom subdictionaries of D. Formally,

$$\mu_s(D) = (\min \{\varsigma_s(D_{\Lambda}) : \Lambda \subseteq [k], |\Lambda| = s\})^2,$$

where $\varsigma_s(A)$ is the s th singular value of A.

The s-incoherence can used to guarantee that sparse codes are stable in a certain sense.

We now introduce some key parameter-and-data-dependent properties. The first property regards the sparsity of the encoder on a sample $\mathbf{x} = (x_1, \dots, x_m)$.

Definition 2 (s-sparsity) If every point x_i in the set of points \mathbf{x} satisfies $\|\varphi_D(x_i)\|_0 \leq s$, then φ_D is s-sparse on \mathbf{x} . More concisely, the boolean expression s-sparse($\varphi_D(\mathbf{x})$) is true.

³Empirically we have observed that stochastic gradient approaches like the one in Mairal et al. (2012) perform very well.

This property is critical as the learning bounds will exploit the observed sparsity level over the training sample. The following collection of properties also will be useful.

Definition 3 (s-margin) Given a dictionary D and a point $x_i \in B_{\mathbb{R}^d}$, the s-margin of D on x_i is

$$\operatorname{margin}_{s}(D, x_{i}) := \max_{\substack{\mathcal{I} \subseteq [k] \\ |\mathcal{I}| = k - s}} \min_{j \in \mathcal{I}} \left\{ \lambda - \left| \langle D_{j}, x_{i} - D\varphi_{D}(x_{i}) \rangle \right| \right\}.$$

The sample version of the s-margin is the maximum s-margin that holds for all points in \mathbf{x} , or the s-margin of D on \mathbf{x} :

$$\operatorname{margin}_s(D, \mathbf{x}) := \min_{x_i \in \mathbf{x}} \operatorname{margin}_s(D, x_i).$$

The importance of these s-margin properties flows directly from the upcoming Sparse Coding Stability Theorem (Theorem 1). Intuitively, if the s-margin of D on x is high, then there is a set of (k-s) inactive atoms that are poorly correlated with the optimal residual $x - D\varphi_D(x)$, and hence these atoms are far from being included in the set of active atoms. More formally, margin_s (D, x_i) is equal to the (s+1)th smallest element of the set of k elements $\{\lambda - |\langle D_j, x_i - D\varphi_D(x_i)\rangle|\}_{j \in [k]}$. Note that if $\|\varphi_D(x_i)\|_0 = s$, we can use the $(s+\rho)$ -margin for any integer $\rho \geq 0$. Indeed, $\rho > 0$ is justified when $\varphi_D(x_i)$ has only s non-zero dimensions but for precisely one index j^* outside the support set $|\langle D_{j^*}, x_i - D\varphi_D(x_i)\rangle|$ is arbitrarily close to λ . In this scenario, the s-margin of D on x_i is trivially small; however, the (s+1)-margin is non-trivial because the max in the definition of the margin will remove j^* from the min's choices \mathcal{I} . Empirical evidence shown in Section 6 suggests that even when ρ is small, the $(s+\rho)$ -margin is not too small.

Sparse coding stability The first result of this work is a fundamental stability result for the LASSO. In addition to being critical in motivating the presented conditions, the result may be of interest in its own right.

Theorem 1 (Sparse Coding Stability) Let dictionaries $D, \tilde{D} \in \mathcal{D}$ satisfy $\mu_s(D), \mu_s(\tilde{D}) \geq \mu$ and $\|D - \tilde{D}\|_2 \leq \varepsilon$ for some $\mu > 0$, and let $x \in B_{\mathbb{R}^d}$. Suppose that there exists an index set $\mathcal{I} \subseteq [k]$ of k - s indices such that for all $i \in \mathcal{I}$:

$$|\langle D_i, x - D\varphi_D(x)\rangle| < \lambda - \tau$$
 (5)

for

$$\varepsilon \le \frac{\tau^2 \lambda}{43} \quad . \tag{6}$$

Then the following stability bound holds:

$$\|\varphi_D(x) - \varphi_{\tilde{D}}(x)\|_2 \le \frac{3\varepsilon\sqrt{s}}{\lambda u}$$
.

Furthermore, if $\varepsilon = \frac{{\tau'}^2 \lambda}{43}$ for $\tau' < \tau$, then for all $i \in \mathcal{I}$:

$$\left| \langle \tilde{D}_i, x - \tilde{D} \varphi_{\tilde{D}}(x) \rangle \right| \leq \lambda - (\tau - \tau').$$

Thus, some margin, and hence sparsity, is retained after perturbation.

Condition (5) means that at least k-s inactive atoms in the coding $\varphi_D(x)$ do not have too high absolute correlation with the residual $x-D\varphi_D(x)$. We refer to the right-hand side of (6) as the permissible radius of perturbation (PRP) because it indicates the maximum amount of perturbation for which the theorem can guarantee encoder stability. In short, the theorem says that if problem (1) admits a stable sparse solution, then a small perturbation to the dictionary will not change the fact that a certain set of k-s atoms remains inactive in the new solution. The theorem further states that the perturbation to the solution will be bounded by a constant factor times the size of the perturbation, where the constant depends on the s-incoherence, the amount of ℓ_1 -regularization, and the sparsity level.

The proof of Theorem 1 is quite long, and so we leave all but the following high-level sketch to Appendix A.

Proof sketch: First, we show that the solution $\varphi_{\tilde{D}}(x)$ is s-sparse and, in particular, has support contained in the complement of \mathcal{I} . Second, we reframe the LASSO as a quadratic program (QP). By exploiting the convexity of the QP and the fact that both solutions have their support contained in a set of s atoms, simple linear algebra yields the desired stability bound. In our view, the first step is much more difficult than the second. Our strategy for the first step has four planks:

- (1) OPTIMAL VALUE STABILITY: The two problems' optimal objective values are close; this is an easy consequence of the closeness of D and \tilde{D} .
- (2) STABILITY OF NORM OF RECONSTRUCTOR: The *norms* of the optimal reconstructors $(D\varphi_D(x))$ and $\tilde{D}\varphi_{\tilde{D}}(x)$ of the two problems are close. We show this using OPTIMAL VALUE STABILITY and

$$(x - D\varphi_D(x))^T D\varphi_D(x) = \lambda \|\varphi_D(x)\|_1, \tag{7}$$

the latter of which can be shown via convex duality Osborne et al. (2000).

- (3) RECONSTRUCTOR STABILITY: The optimal reconstructors of the two problems are close. This fact can be shown to be a consequence of STABILITY OF NORM OF RECONSTRUCTOR, using the ℓ_1 norm's convexity and the equality (7).
- (4) PRESERVATION OF SPARSITY: The solution to the perturbed problem also is supported on the complement of \mathcal{I} . To show this, it is sufficient to show that the absolute correlation of each atom \tilde{D}_i ($i \in \mathcal{I}$) with the residual in the perturbed problem is less than λ . This last claim is a relatively easy consequence of RECONSTRUCTOR STABILITY.

Although we do not make use of it in this work, under considerably stronger conditions, we can achieve a similar stability bound with a far smaller PRP.

Theorem 2 (Restricted Stability⁴) Let dictionaries $D, \tilde{D} \in \mathcal{D}$ satisfy $\mu_s(D), \mu_s(\tilde{D}) \geq \mu$ and $\|D - \tilde{D}\|_2 \leq \varepsilon$ for some $\mu > 0$, and let $x \in B_{\mathbb{R}^d}$. Suppose that there exist $\tau > 0$ and $s \in [k]$ satisfying:

(ii)
$$|(\varphi_D(x))_j| > \tau \text{ for all } j \in \text{supp } \varphi_D(x), \tag{9}$$

(iii)
$$|\langle D_j, x - D\varphi_D(x)\rangle| < \lambda - \tau \text{ for all } j \notin \operatorname{supp} \varphi_D(x); \tag{10}$$

$$\varepsilon \le \frac{\tau \mu}{\frac{s+\mu}{\lambda} + \sqrt{s} + \mu} \quad . \tag{11}$$

Then the supports of the optimal solutions are identical:

$$\operatorname{supp}(\varphi_D(x)) = \operatorname{supp}(\varphi_{\tilde{D}}(x)) ,$$

and the following stability bound holds:

$$\|\varphi_D(x) - \varphi_{\tilde{D}}(x)\|_2 \le \frac{\varepsilon}{\mu} \left(\frac{\sqrt{s}}{\lambda} + 1\right).$$

Theorem 2 applies provided that every non-zero coefficient of $\varphi_D(x)$ has magnitude bounded above zero and every unused atom is far from being brought into the optimal solution in the sense of (10) (note this inequality is identical to (5). The conditions of Theorem 2 are more demanding than those of Theorem 1 because in the former, there is exactly one choice of the inactive set \mathcal{I} , namely, $[k] \setminus \text{supp } \varphi_D(x)$. If this choice does not yield sufficient margin, we are out of luck. Nevertheless, when the conditions of the Theorem 2 do hold, the PRP's dependence on the margin-like property τ is considerably reduced from quadratic in Theorem 1 to only linear in Theorem 2. Although the PRP now depends on s and μ , both of these properties often are well-behaved (a small s is desired and typically leads to a large μ), whereas τ is a wild-card on which minimum dependence is desired.

Proof sketch: Our strategy is to show that there is a unique solution to the perturbed problem, defined in terms of the optimality conditions of the LASSO (see conditions L1 and L2 of (Asif and Romberg, 2010)), and this solution has the same support as the solution to the original problem. As a result, the perturbed solution's proximity to the original solution is governed in part by a condition number μ of a linear system of s variables.

2.1 Main results

The following notation will aid and abet the below results and the subsequent analysis. Recall that the loss l is bounded by b and L-Lipschitz in its second argument. Also recall that \mathcal{F} is the set of predictive sparse coding hypothesis functions $f(x) = \langle w, \varphi_D(x) \rangle$ indexed by $D \in \mathcal{D}$ and $w \in \mathcal{W}$. For $f \in \mathcal{F}$, define $l(\cdot, f) : \mathcal{Y} \times \mathbb{R}^d \to [0, b]$ as the loss-composed function $(y, x) \mapsto l(y, f(x))$. Let $l \circ \mathcal{F}$ be the class of such functions induced by the choice of \mathcal{F} and l. A probability measure P operates on functions and loss-composed functions as:

$$Pf = \mathsf{E}_{(x,y) \sim P} f(x) \qquad \qquad Pl(\cdot,f) = \mathsf{E}_{(x,y) \sim P} l(y,f(x)).$$

Similarly, an empirical measure $P_{\mathbf{z}}$ associated with sample \mathbf{z} operates on functions and loss-composed functions as:

$$P_{\mathbf{z}}f = \frac{1}{m} \sum_{i=1}^{m} f(x_i)$$
 $P_{\mathbf{z}}l(\cdot, f) = \frac{1}{m} \sum_{i=1}^{m} l(y_i, f(x_i)).$

Finally, when provided a training sample \mathbf{z} , the hypothesis returned by the learner will be referred to as $\hat{f}_{\mathbf{z}}$. Note that $\hat{f}_{\mathbf{z}}$ is random, but $\hat{f}_{\mathbf{z}}$ becomes a fixed function upon conditioning on \mathbf{z} .

⁴In the previous version of this paper (Mehta and Gray, 2012), the Theorem 1 of the current paper did not exist and the Theorem 2 of this paper was labeled as Theorem 1.

Classically speaking, the overcomplete setting is the *modus operandi* in sparse coding. In this setting, an overcomplete basis is learned which will be used parsimoniously in coding individual points. The next result bounds the generalization error in the overcomplete setting. The $\tilde{O}(\cdot)$ notation hides $\log(\log(\cdot))$ terms and assumes that $r \leq m^{\min\{d,k\}}$.

Theorem 3 (Overcomplete Learning Bound) With probability at least $1 - \delta$ over $\mathbf{z} \sim P^m$, for any $s \in [k]$ and any $f = (D, w) \in \mathcal{F}$ satisfying s-sparse $(\varphi_D(\mathbf{x}))$ and

$$m > \frac{387}{\text{margin}_s(D, \mathbf{x})^2 \lambda},$$

the generalization error $(P - P_z)l(\cdot, f)$ is

$$\tilde{O}\left(b\sqrt{\frac{dk\log m + \log\frac{1}{\delta}}{m}} + \frac{b}{m}\left(dk\log\frac{1}{\mathrm{margin}_s^2(D,\mathbf{x})\cdot\lambda}\right) + \frac{L}{m}\left(\frac{r\sqrt{s}}{\lambda\mu_s(D)}\right)\right). \tag{12}$$

Note that this bound also applies to the particular hypothesis $\hat{f}_{\mathbf{z}} = (\hat{D}_{\mathbf{z}}, \hat{w}_{\mathbf{z}})$ learned from the training sample.

Often in learning problems, we first map the data implicitly to a space of very high dimension or even infinite dimension and use kernels for efficient computations. In these cases where $d \gg k$ or d is infinite, it is unacceptable for any learning bound to exhibit dependence on d. It is possible to untether the analysis from d by using the s-margin of the learned dictionary $\hat{D}_{\mathbf{z}}$ on a second, unlabeled sample. In the infinite-dimensional setting, the following dimension-free learning bound holds.

Theorem 4 (Infinite-Dimensional Learning Bound) With probability at least $1 - \delta$ over a labeled m-sample $\mathbf{z} \sim P^m$ and a second, unlabeled sample $\mathbf{x}'' \sim \Pi^m$, if an algorithm learns hypothesis $\hat{f}_{\mathbf{z}} = (\hat{D}_{\mathbf{z}}, \hat{w}_{\mathbf{z}})$ such that $\varphi_{\hat{D}_{\mathbf{z}}}$ is s-sparse on $(\mathbf{x} \cup \mathbf{x}'')$, $\mu_{2s}(\hat{f}_{\mathbf{z}}) > 0$, and

$$m \ge \frac{43}{\mathrm{margin}_s^2(\hat{D}_{\mathbf{z}}, \mathbf{x} \cup \mathbf{x''}) \cdot \lambda},$$

then the generalization error $(P - P_z)l(\cdot, \hat{f}_z)$ is

$$\tilde{O}\left(\frac{L}{\sqrt{m}}\left(\frac{rk\sqrt{s}}{\mu_{2s}(\hat{f}_{\mathbf{z}})}\right) + b\sqrt{\frac{(k^2 + \log\frac{1}{\delta})\log m}{m}} + \frac{L}{m}\left(\frac{r\sqrt{s}}{\lambda\mu_s(\hat{f}_{\mathbf{z}})}\right)\right). \tag{13}$$

2.2 Discussion of Theorems 3 and 4

The results highlight the central role of the stability of the sparse encoder. The presented bounds are data-dependent and exploit properties relating to the training sample and the learned hypothesis. Since $k \geq d$ in the overcomplete setting, an ideal learning bound has minimal dependence on k. The $\frac{1}{m}$ term of the learning bound for the overcomplete setting (12) exhibits square root dependence on both the size of the dictionary k and the ambient dimension d. It is unclear whether further improvement is possible, even in the reconstructive setting. The two known results in the reconstructive setting were established first by Maurer and Pontil (2010) and later by Vainsencher

et al. (2011), as mentioned in the Introduction. The infinite-dimensional setting learning bound (13) is dimension free, with linear dependence on k, square root dependence on s, and inverse dependence on the 2s-incoherence $\mu_{2s}(\hat{f}_{\mathbf{z}})$. While both bounds exhibit dependence on the sparsity level s, the sparsity level appears to be much more significant in the infinite-dimensional setting.

Let us compare these bounds to the reconstructive setting, starting with the overcomplete regime. The first term of (12) matches the slower of the rates shown by Vainsencher et al. (2011) for the unsupervised case. Vainsencher et al. also showed fast rates of $\frac{dk}{m}$ (plus a small fraction of the observed empirical risk), but in the predictive setting it is an open question whether similar fast rates are possible. The second term of (12) represents the error in approximating the estimator via an $(\varepsilon = \frac{1}{m})$ -cover of the space of dictionaries. This term reflects the stability of the sparse codes with respect to dictionary perturbations, as quantified by the Sparse Coding Stability Theorem (Theorem 1). The reason for the lower bound on m is that the ε -net used to approximate the space of dictionaries needs to be fine enough to satisfy the PRP condition (6) of the Sparse Coding Stability Theorem. Hence, both this lower bound and the second term are determined primarily by the Sparse Coding Stability Theorem, and so with this proof strategy the extent to which the Sparse Coding Stability Theorem cannot be improved also indicates the extent to which Theorem 3 cannot be improved.

Shifting to the infinite-dimensional setting, Maurer and Pontil (2010) previously showed the generalization bound (2) for unsupervised (ℓ_1 -regularized) sparse coding. Comparing their result to (13) and neglecting regularization parameters, the dimension-free bound in the predictive case is larger by a factor of $\frac{\sqrt{s}}{\mu_{2s}(\hat{f}_z)}$. It is unclear whether either of the terms in this factor are avoidable in the predictive setting. At least from our analysis, it appears that the $\frac{\sqrt{s}}{\mu_{2s}(\hat{f}_z)}$ factor is the price one pays for encoder stability. Critically, encoder stability is not necessary in the reconstructive setting because stability in loss (reconstruction error) requires only stability in the norm of the residual to the LASSO problem rather than stability in the value of the solution to the problem. Stability of the norm of the residual is readily obtainable without any of the incoherence, sparsity, and margin conditions used here.

Remarks on conditions One may wonder about typical values for the various hypothesis-and-data-dependent properties in Theorems 3 and 4. In practical applications of reconstructive and predictive sparse coding, the regularization parameter λ is set to ensure that s is small relative to the dimension d. As a result, both incoherences $\mu_s(D)$ and $\mu_{2s}(D)$ for the learned dictionary can be expected to be bounded away from zero. A sufficiently large s-incoherence certainly is necessary if one hopes for any amount of stability of the class of sparse coders with respect to dictionary perturbations. Since our path to reaching Theorems 3 and 4 passes through the Sparse Coding Stability Theorem (Theorem 1), it seems that a drastically different strategy needs to be used if it is possible to avoid dependence on $\mu_s(D)$ in the learning bounds.

A curious aspect of both learning bounds is their dependence on the s-margin term margin_s(D, \mathbf{x}). Suppose that a dictionary is learned which is s-sparse on the training sample \mathbf{x} , and s is the lowest such integer for which this holds. It may not always be the case that the s-margin is bounded away from zero because for some points a small collection of inactive atoms may be very close to being brought into the optimal solution (the code); however, we can instead use the $(s + \rho)$ -margin for some small positive integer ρ for which the $(s + \rho)$ -margin is non-trivial. In Section 6, we gather empirical evidence that such a non-trivial $(s + \rho)$ -margin does exist, for small ρ , when learning

predictive sparse codes on real data. Hence, there is evidence that predictive sparse coding learns a dictionary with high s-incoherence $\mu_s(D)$ and non-trivial s-margin margin_s (D, \mathbf{x}) on the training sample, for low s.

If one entertains a mixture of ℓ_1 and ℓ_2 norm regularization, $\lambda_1|\cdot||_1 + \frac{1}{2}\lambda_2||\cdot||_2^2$, as in the elastic net (Zou and Hastie, 2005), fall-back guarantees are possible in both scenarios. For small values of λ_2 , this regularizer induces true sparsity similar to the ℓ_1 regularizer. A considerably simpler, data-independent analysis is possible in the overcomplete setting with a final bound that essentially just trades $\mu_s(D)$ for the ℓ_2 norm regularization parameter λ_2 . In the infinite-dimensional setting, a simpler non-data-dependent analysis using our approach would only attain a bound of the larger order $\frac{k^{3/2}}{\lambda_2\sqrt{m}}$.

3 Tools

As before, let **z** be a labeled sample of m points (an m-sample) drawn iid from P. In addition, let \mathbf{z}' be a second labeled m-sample drawn iid from P. In the infinite-dimensional setting, we also will make use of an unlabeled m-sample \mathbf{x}'' drawn iid from the marginal Π . Also, an *epsilon-cover* will be used to refer to the *concept* of an ε -cover but not any specific cover. All epsilon-covers of spaces of dictionaries use the metric induced by the operator norm $\|\cdot\|_2$.

3.1 Symmetrization by ghost sample for random subclasses

The next result is essentially due to Mendelson and Philips (2004); it applies symmetrization by a ghost sample for random subclasses. Our main departure is that we allow the random subclass to depend on a second, unlabeled sample \mathbf{x}'' .

Lemma 1 (Symmetrization by Ghost Sample) Let $\mathcal{F}(\mathbf{z}, \mathbf{x}'') \subset \mathcal{F}$ be a random subclass which can depend on both a labeled sample \mathbf{z} and an unlabeled sample \mathbf{x}'' . Recall that \mathbf{z}' is a ghost sample of m points. If $m \geq \left(\frac{b}{t}\right)^2$, then

$$\Pr_{\mathbf{z},\mathbf{x}''} \left\{ \exists f \in \mathcal{F}(\mathbf{z},\mathbf{x}''), \ (P - P_{\mathbf{z}})l(\cdot,f) \ge t \right\}$$

$$\leq 2\Pr_{\mathbf{z},\mathbf{z}',\mathbf{x}''} \left\{ \exists f \in \mathcal{F}(\mathbf{z},\mathbf{x}''), \ (P_{\mathbf{z}'} - P_{\mathbf{z}})l(\cdot,f) \ge \frac{t}{2} \right\}.$$

For completeness, this lemma is proved in Appendix C. This symmetrization lemma will be applied in both the overcomplete and infinite-dimensional settings to shift the analysis from large deviations of the empirical risk from the expected risk to large deviations of two independent empirical risks: in the overcomplete setting the lemma will be specialized as Proposition 1, and in the infinite-dimensional setting the lemma will be adapted to Proposition 2.

3.2 Rademacher and Gaussian averages and related results

Let $\sigma_1, \ldots, \sigma_m$ be independent Rademacher random variables distributed uniformly on $\{-1, 1\}$, and let $\gamma_1, \ldots, \gamma_m$ be independent Gaussian random variables distributed as $\mathcal{N}(0, 1)$. Denote the collections by $\boldsymbol{\sigma} = (\sigma_1, \ldots, \sigma_m)$ and $\boldsymbol{\gamma} = (\gamma_1, \ldots, \gamma_m)$. Given a sample of m points \mathbf{x} , define the

conditional Rademacher and Gaussian averages of a function class as

$$\mathcal{R}_{m|\mathbf{x}}(\mathcal{F}) = \frac{2}{m} \mathsf{E}_{\boldsymbol{\sigma}} \sup_{f \in \mathcal{F}} \sum_{i=1}^{m} \sigma_{i} f(x_{i}) \quad \text{and} \quad \mathcal{G}_{m|\mathbf{x}}(\mathcal{F}) = \frac{2}{m} \mathsf{E}_{\boldsymbol{\gamma}} \sup_{f \in \mathcal{F}} \sum_{i=1}^{m} \gamma_{i} f(x_{i}).$$

respectively.

Lemmas 2 and 3 below are used near the end of the proof of Theorem 7 of the infinite-dimensional setting, when shifting the analysis from the Gaussian complexity of a loss-composed function class to the Rademacher complexity of the original function class. From Meir and Zhang (2003, Theorem 7), the loss-composed conditional Rademacher average of a function class \mathcal{F} is bounded by the scaled conditional Rademacher average:

Lemma 2 (Rademacher Loss Comparison Lemma) For every function class \mathcal{F} , m-sample \mathbf{x} , and l which is L-Lipschitz continuous in its second argument:

$$\mathcal{R}_{m|\mathbf{z}}(l \circ \mathcal{F}) \leq L\mathcal{R}_{m|\mathbf{x}}(\mathcal{F}).$$

Additionally, from Ledoux and Talagrand (1991, a brief argument following Lemma 4.5), the conditional Rademacher average of a function class \mathcal{F} is bounded up to a constant by the conditional Gaussian average of \mathcal{F} :

Lemma 3 (Rademacher-Gaussian Average Comparison Lemma) For every function class \mathcal{F} and sample of m points \mathbf{x} :

$$\mathcal{R}_{m|\mathbf{x}}(\mathcal{F}) \leq \sqrt{\frac{\pi}{2}} \mathcal{G}_{m|\mathbf{x}}(\mathcal{F}).$$

The next relation is due to Slepian (1962):

Lemma 4 (Slepian's Lemma) Let Ω and Γ be mean zero, separable Gaussian processes⁵ indexed by a set T such that $\mathsf{E}(\Omega_{t_1} - \Omega_{t_2})^2 \leq \mathsf{E}(\Gamma_{t_1} - \Gamma_{t_2})^2$ for all $t_1, t_2 \in T$. Then $\mathsf{E}\sup_{t \in T} \Omega_t \leq \mathsf{E}\sup_{t \in T} \Gamma_t$.

Slepian's Lemma essentially says that if the variance of one Gaussian process is bounded by the variance of another, then the expected maximum of the first is bounded by the expected maximum of the second. This lemma will be used in the proof of Theorem 6 to bound the Gaussian complexity of an analytically difficult function class via a bound on the Gaussian complexity of a related but analytically easier function class.

We also will make use of the following bounded differences inequality due to McDiarmid (1989), in order to shift the analysis in the proof of Theorem 7 to the Rademacher complexity of a certain function class:

Theorem 5 (McDiarmid's Inequality) Let X_1, \ldots, X_m be random variables drawn iid according to a probability measure μ over a space \mathcal{X} . Suppose that a function $f: \mathcal{X}^m \to \mathbb{R}$ satisfies

$$\sup_{x_1, \dots, x_m, x_i' \in \mathcal{X}} |f(x_1, \dots, x_n) - f(x_1, \dots, x_{i-1}, x_i', x_{i+1}, \dots, x_m)| \le c_i$$

for any $i \in [m]$. Then

$$\Pr_{X_1,...,X_n} \{ f(X_1,...,X_n) - \mathsf{E} f(X_1,...,X_n) \ge t \} \le \exp\left(-2t^2 / \sum_{i=1}^m c_i^2\right).$$

 $^{{}^{5}\{\}Omega_{t}\}_{t\in T}$ is a Gaussian process with index set T if the collection is jointly Gaussian in the sense that every finite linear combination of the variables is Gaussian.

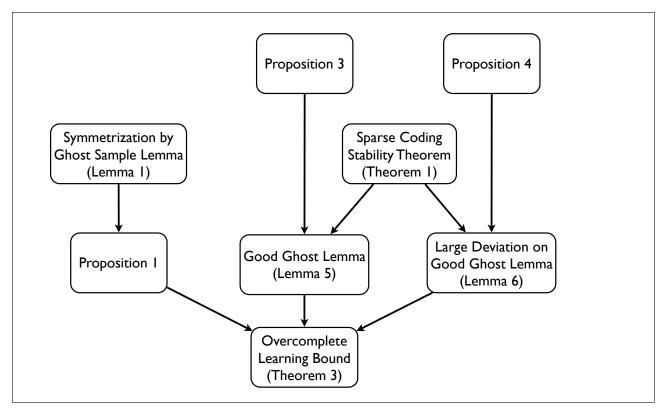


Figure 1: Proof flowchart for the Overcomplete Learning Bound (Theorem 3).

4 Overcomplete setting

The overcomplete setting is classically the more popular regime, and in this setting useful learning bounds may exhibit dependence on both the dimension d and the dictionary size k. At a high level, our strategy for the overcomplete case learning bound is to construct an epsilon-cover over a subclass of the space of functions $\mathcal{F} := \{f = (D, w) : D \in \mathcal{D}, w \in \mathcal{W}\}$ and to show that the metric entropy of this subclass is of order dk. The main difficulty is that an epsilon-cover over \mathcal{D} need not approximate \mathcal{F} to any degree, unless one has a notion of encoder stability. Our analysis effectively will be concerned only with a training sample and a ghost sample, and it is similar in style to the luckiness framework of Shawe-Taylor et al. (1998). If we observe that the sufficient conditions for encoder stability hold true on the training sample, then it is enough to guarantee that most points in a ghost sample also satisfy these conditions (at a weaker level). Figure 1 exhibits the high-level flow of the proof of Theorem 3.

4.1 Useful conditions and subclasses

Let $\tilde{\mathbf{x}} \subseteq_{\eta} \mathbf{x}$ indicate that $\tilde{\mathbf{x}}$ is a subset of \mathbf{x} with at most η elements of \mathbf{x} removed. This notation is identical to Shawe-Taylor et al. (1998)'s notation from the luckiness framework.

Our bounds will require a crucial PRP-based condition that depends on both the learned dic-

tionary and the training sample:

$$\mathrm{margin}_s(D,\mathbf{x}) \geq \iota(\lambda,\mu,\varepsilon) \qquad \qquad \text{for } \iota(\lambda,\mu,\varepsilon) = \sqrt{\frac{387\varepsilon}{\lambda}}.$$

For brevity we will refer to ι with its parameters implicit; the dependence on ε , λ , and μ will not be an issue because we first develop bounds with these quantities fixed a priori. Lastly, for $\mu > 0$ define $\mathcal{D}_{\mu} := \{D \in \mathcal{D} : \mu_s(D) \geq \mu\}$ and $\mathcal{F}_{\mu} := \{f = (D, w) \in \mathcal{F} : D \in \mathcal{D}_{\mu}\}.$

4.2 Learning bound

The following proposition is simply a specialization of Lemma 1 with \mathbf{x}'' taken as the empty set and $\mathcal{F}(\mathbf{z}, \mathbf{x}'') := \{ f \in \mathcal{F}_{\mu} : \lceil \max_{s}(D, \mathbf{x}) > \iota \rceil \}.$

Proposition 1 If $m \geq \left(\frac{b}{t}\right)^2$, then

$$\begin{split} & \Pr_{\mathbf{z}} \left\{ \exists f \in \mathcal{F}_{\mu}, \ \left[\operatorname{margin}_{s}(D, \mathbf{x}) \right) > \iota \right] \ \mathbf{and} \ \left((P - P_{\mathbf{z}}) l(\cdot, f) > t \right) \right\} \\ & \leq 2 \operatorname{Pr}_{\mathbf{z} \, \mathbf{z}'} \left\{ \exists f \in \mathcal{F}_{\mu}, \ \left[\operatorname{margin}_{s}(D, \mathbf{x}) > \iota \right] \ \mathbf{and} \ \left((P_{\mathbf{z}'} - P_{\mathbf{z}}) l(\cdot, f) > t / 2 \right) \right\}. \end{split}$$

In the RHS of the above, let the event whose probability is being measured be

$$J := \left\{ \mathbf{z} \, \mathbf{z}' : \exists f \in \mathcal{F}_{\mu}, \, \left[\operatorname{margin}_{s}(D, \mathbf{x}) > \iota \right] \, \operatorname{\mathbf{and}} \, \left(P_{\mathbf{z}'} - P_{\mathbf{z}} \right) l(\cdot, f) > t/2) \right\}.$$

Define Z as the event that there exists a hypothesis with stable codes on the original sample, in the sense of the Sparse Coding Stability Theorem (Theorem 1), but more than $\eta = \eta(m, d, k, D, \mathbf{x}, \delta)$ points⁶ of the ghost sample whose codes are not guaranteed stable by the Sparse Coding Stability Theorem:

$$Z := \left\{ \mathbf{z} \, \mathbf{z}' : \begin{array}{c} \exists f \in \mathcal{F}_{\mu}, \ \left[\operatorname{margin}_{s}(D, \mathbf{x}) > \iota \right] \\ \text{ and } \left(\nexists \, \tilde{\mathbf{x}} \subseteq_{\eta} \, \mathbf{x}' \left[\operatorname{margin}_{s}(D, \tilde{\mathbf{x}}) > \frac{1}{3} \operatorname{margin}_{s}(D, \mathbf{x}) \right] \right) \end{array} \right\}.$$

Our strategy will be to show that Pr(J) is small by use of the fact that

$$\Pr(J) = \Pr(J \cap \bar{Z}) + \Pr(J \cap Z) \le \Pr(J \cap \bar{Z}) + \Pr(Z),$$

a strategy which thus far is similar to the beginning of Shawe-Taylor et al.'s proof of the main luckiness framework learning bound (see Shawe-Taylor et al., 1998, Theorem 5.22). We now show that each of $\Pr(Z)$ and $\Pr(J \cap \bar{Z})$ is small in turn.

The imminent Good Ghost Lemma shadows Shawe-Taylor et al. (1998)'s notion of probable smoothness and provides a bound on Pr(Z).

Lemma 5 (Good Ghost) Fix $\mu, \lambda > 0$ and $s \in [k]$. With probability at least $1 - \delta$ over an m-sample $\mathbf{x} \sim P^m$ and a second m-sample $\mathbf{x}' \sim P^m$, for any $D \in \mathcal{D}_{\mu}$ for which φ_D is s-sparse on \mathbf{x} , at least $m - \eta(m, d, k, D, \mathbf{x}, \delta)$ points $\tilde{\mathbf{x}} \subseteq \mathbf{x}'$ satisfy $[\text{margin}_s(D, \tilde{\mathbf{x}}) > \frac{1}{3} \text{margin}_s(D, \mathbf{x})]$, for

$$\eta(m, d, k, D, \mathbf{x}, \delta) := dk \log \frac{3096}{\mathrm{margin}_s^2(D, \mathbf{x}) \cdot \lambda} + \log(2m + 1) + \log \frac{1}{\delta}.$$

⁶We use the shorthand $\eta = \eta(m, d, k, D, \mathbf{x}, \delta)$ for conciseness.

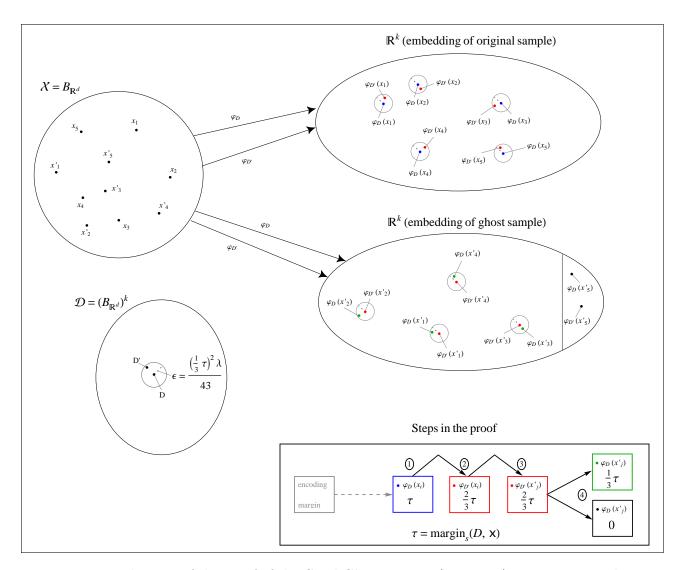


Figure 2: Visualization of the proof of the Good Ghost Lemma (Lemma 5). Best seen in color.

Proof: Figure 2 illustrates the proof. By the assumptions of the lemma, consider an arbitrary dictionary D satisfying $\mu_s(D) \geq \mu$ and s-sparse($\varphi_D(\mathbf{x})$). The goal is to guarantee with high probability that all but η points of the ghost sample are coded by φ_D with s-margin of at least $\frac{1}{3} \operatorname{margin}_s(D, \mathbf{x})$.

Let $\varepsilon = \frac{(\frac{1}{3} \operatorname{margin}_s(D, \mathbf{x}))^2 \cdot \lambda}{43}$, and consider a minimum-cardinality ε -proper cover \mathcal{D}' of \mathcal{D}_{μ} . Let D' be a candidate element of \mathcal{D}' satisfying $||D - D'||_2 \le \varepsilon$. Then the Sparse Coding Stability Theorem (Theorem 1) implies that the coding margin of D' on \mathbf{x} retains over two-thirds the coding margin of D on \mathbf{x} ; that is, $[\operatorname{margin}_s(D', \mathbf{x}) > \frac{2}{3} \operatorname{margin}_s(D, \mathbf{x})]$.

Furthermore, we can and will show that most points from the ghost sample satisfy $[\text{margin}_s(D',\cdot)>$

 $\frac{2}{3}$ margin $_s(D, \mathbf{x})$]. Let $\mathcal{F}_D^{\text{marg}} := \{f_{D, \tau}^{\text{marg}} | \tau \in \mathbb{R}_+\}$ be the class of threshold functions defined via

$$f_{D,\tau}^{\text{marg}}(x) := \begin{cases} 1; & \text{if } \text{margin}_s(D,x) > \tau, \\ 0; & \text{otherwise.} \end{cases}$$

Since the VC dimension of the one-dimensional threshold functions is 1, it follows that the VC(\mathcal{F}_D^{marg}) = 1. By using the VC dimension of \mathcal{F}_D^{marg} and the standard permutation argument of Vapnik and Chervonenkis (1968, Proof of Theorem 2), it follows that for a single, fixed element of \mathcal{D}' , with probability at least $1-\delta$ at most $\log(2m+1)+\log\frac{1}{\delta}$ points from a ghost sample will violate the margin inequality in question. Hence, by the bound on the proper covering numbers provided by Proposition 3 (see Appendix F), we can we can guarantee for all candidate members $D' \in \mathcal{D}'$ that with probability $1-\delta$ at most

$$\eta(m, d, k, D, \mathbf{x}, \delta) = dk \log \frac{3096}{\text{margin}_s^2(D, \mathbf{x}) \cdot \lambda} + \log(2m + 1) + \log \frac{1}{\delta}$$

points from the ghost sample violate the s-margin inequality. Thus, for arbitrary $D' \in \mathcal{D}'$ satisfying the conditions of the lemma, with probability $1 - \delta$ at most $\eta(m, d, k, D, \mathbf{x}, \delta)$ points from the ghost sample violate $\left[\operatorname{margin}_s(D', \cdot) > \frac{2}{3} \operatorname{margin}_s(D, \mathbf{x}) \right]$.

Finally, consider the at least $m - \eta(m,d,k,D,\mathbf{x},\delta)$ points in the ghost sample that satisfy $\left[\operatorname{margin}_s(D',\cdot) > \frac{2}{3}\operatorname{margin}_s(D,\mathbf{x})\right]$. Since $\|D'-D\|_2 \leq \frac{(\frac{1}{3}\operatorname{margin}_s(D,\mathbf{x}))^2 \cdot \lambda}{43}$, the Sparse Coding Stability Theorem (Theorem 1) implies that these points satisfy $\left[\operatorname{margin}_s(D,\cdot) > \frac{1}{3}\operatorname{margin}_s(D,\mathbf{x})\right]$.

It remains to bound $\Pr(J \cap \bar{Z})$.

Lemma 6 (Large Deviation on Good Ghost) Define $\varpi := t/2 - \left(2L\beta + \frac{b\eta}{m}\right)$ and $\beta := \frac{\varepsilon}{\lambda} \left(1 + \frac{3r\sqrt{s}}{\mu}\right)$. Then

$$\Pr(J \cap \bar{Z}) \le \left(\frac{8(r/2)^{1/(d+1)}}{\varepsilon}\right)^{(d+1)k} \exp(-m\varpi^2/(2b^2)).$$

Equivalently, the difference between the loss on \mathbf{z} and the loss on \mathbf{z}' is greater than $\varpi + 2L\beta + \frac{b\eta}{m}$ with probability at most $\left(\frac{8(r/2)^{1/(d+1)}}{\varepsilon}\right)^{(d+1)k} \exp(-m\varpi^2/(2b^2))$.

Proof: Let GOOD GHOST represent the condition that all but at most η points in the ghost sample \mathbf{x}' satisfy $\left[\operatorname{margin}_s(D,\cdot) > \frac{1}{3}\operatorname{margin}_s(D,\mathbf{x})\right]$ (and hence are "good"), and let BAD GHOST be true if and only if GOOD GHOST is false.

First, note that the event $J \cap \bar{Z}$ is a subset of the event

To see this, suppose that $\mathbf{z}\mathbf{z}' \in J \cap \bar{Z}$. Since $\mathbf{z}\mathbf{z}' \in J$, there is a particular f = (D, w) in \mathcal{F}_{μ} satisfying $[\operatorname{margin}_{s}(D, \mathbf{x}) > \iota]$ and $(P_{\mathbf{z}'} - P_{\mathbf{z}})l(\cdot, f) > t/2$. Since $\mathbf{z}\mathbf{z}' \in \bar{Z}$, there is no function in \mathcal{F}_{μ} satisfying both $[\operatorname{margin}_{s}(D, \mathbf{x}) > \iota]$ and BAD GHOST. But since the particular f satisfies $[\operatorname{margin}_{s}(D, \mathbf{x}) > \iota]$, it cannot also satisfy BAD GHOST; thus, the particular f satisfies GOOD GHOST.

Bounding the probability of the event R is equivalent to bounding the probability of a large deviation (i.e. $((P_{\mathbf{z}'} - P_{\mathbf{z}})l(\cdot, f) > t/2))$ for the random subclass:

$$\tilde{\mathcal{F}}(\mathbf{x}, \mathbf{x}') := \left\{ \begin{array}{l} f \in \mathcal{F}_{\mu} : \left[\operatorname{margin}_{s}(D, \mathbf{x}) > \iota \right] \text{ and } \\ \left(\exists \tilde{\mathbf{x}} \subseteq_{\eta} \mathbf{x}', \left[\operatorname{margin}_{s}(D, \tilde{\mathbf{x}}) > \frac{1}{3} \operatorname{margin}_{s}(D, \mathbf{x}) \right] \right) \end{array} \right\}.$$

Let $\mathcal{F}_{\varepsilon} = \mathcal{D}_{\varepsilon} \times \mathcal{W}_{\varepsilon}$, where $\mathcal{D}_{\varepsilon}$ is a minimum-cardinality proper ε -cover of \mathcal{D}_{μ} and $\mathcal{W}_{\varepsilon}$ is a minimum-cardinality ε -cover of \mathcal{W} . It is sufficient to bound the probability of a large deviation for all of $\mathcal{F}_{\varepsilon}$ and to then consider the maximum difference between an element of $\tilde{\mathcal{F}}(\mathbf{x}, \mathbf{x}')$ and its closest representative in $\mathcal{F}_{\varepsilon}$. Clearly, for each $f = (D, w) \in \tilde{\mathcal{F}}(\mathbf{x}, \mathbf{x}')$, there is a $f' = (D', w') \in \mathcal{F}_{\varepsilon}$ satisfying $\|D - D'\|_2 \le \varepsilon$ and $\|w - w'\|_2 \le \varepsilon$. If ε is sufficiently small, then for all but η of the points x_i in the ghost sample (and for all points x_i of the original sample) it is guaranteed that

$$\begin{aligned} |\langle w, \varphi_D(x_i) \rangle - \langle w', \varphi_{D'}(x_i) \rangle| &\leq |\langle w - w', \varphi_D(x_i) \rangle| + |\langle w', \varphi_D(x_i) - \varphi_{D'}(x_i) \rangle| \\ &\leq \frac{\varepsilon}{\lambda} + r \frac{3\varepsilon\sqrt{s}}{\lambda\mu} \\ &= \frac{\varepsilon}{\lambda} \left(1 + \frac{3r\sqrt{s}}{\mu} \right) = \beta, \end{aligned}$$

where the second inequality follows from the Sparse Coding Stability Theorem (Theorem 1). Trivially, for the rest of the points x_i in the ghost sample there is a coarse guarantee that $\|\varphi_D(x_i) - \varphi_{D'}(x_i)\|_2 \leq \frac{2}{\lambda}$. Hence, on the original sample:

$$\frac{1}{m} \sum_{i=1}^{m} \left| l(y_i, \langle w, \varphi_D(x_i) \rangle) - l(y_i, \langle w', \varphi_{D'}(x_i) \rangle) \right| \le L\beta,$$

and on the ghost sample:

$$\frac{1}{m} \sum_{i=1}^{m} \left| l(y_i', \langle w, \varphi_D(x_i') \rangle) - l(y_i', \langle w', \varphi_{D'}(x_i') \rangle) \right| \\
\leq \frac{L}{m} \sum_{i \text{ GOOD}} \left| \langle w, \varphi_D(x_i) \rangle - \langle w', \varphi_{D'}(x_i) \rangle \right| + \frac{1}{m} \sum_{i \text{ BAD}} \left| l(y_i', \langle w, \varphi_D(x_i') \rangle) - l(y_i', \langle w', \varphi_{D'}(x_i') \rangle) \right| \\
\leq L\beta + \frac{b\eta}{m},$$

where GOOD denotes the at least $m - \eta$ points of the ghost sample for which the Sparse Coding Stability Theorem (Theorem 1) applies, and BAD denotes the complement thereof. To conclude the above argument, the difference between the losses of f and f' on the double sample will be at most $2L\beta + \frac{b\eta}{m}$.

Now, if ν is the absolute deviation between the loss of f on the original sample versus its loss on the ghost sample, then the absolute deviation between the loss of f' = (D', w') on the original sample and the loss of f' on the ghost sample must be at least

$$\nu - \left(2L\beta + \frac{b\eta}{m}\right).$$

Consequently, if $\nu > t/2$, then the absolute deviation between the loss of f' on the original sample and the loss of f' on the ghost sample must be at least $t/2 - \left(2L\beta + \frac{b\eta}{m}\right)$. To bound the probability of R it therefore is sufficient to control

$$\Pr_{\mathbf{z}\,\mathbf{z}'}\left\{\exists f=(D',w')\in\mathcal{D}_{\varepsilon}\times\mathcal{W}_{\varepsilon},\ (P_{\mathbf{z}'}-P_{\mathbf{z}})l(\cdot,f)>t/2-\left(2L\beta+\frac{b\eta}{m}\right)\right\}.$$

We first handle the case of a fixed $f = (D', w') \in \mathcal{D}_{\varepsilon} \times \mathcal{W}_{\varepsilon}$. Applying Hoeffding's inequality to the random variable $l(y_i, f(x_i)) - l(y_i', f(x_i'))$, with range in [-b, b], yields:

$$\Pr_{\mathbf{z}\mathbf{z}'}\left\{ (P_{\mathbf{z}'} - P_{\mathbf{z}})l(\cdot, f) > \varpi \right\} \le \exp(-m\varpi^2/(2b^2)),$$

for $\varpi := t/2 - \left(2L\beta + \frac{b\eta}{m}\right)$. By way of a proper covering number bound of $\mathcal{D}_{\varepsilon} \times \mathcal{W}_{\varepsilon}$ (see Proposition 4) and the union bound, this result can be extended over all of $\mathcal{D}_{\varepsilon} \times \mathcal{W}_{\varepsilon}$:

$$\Pr_{\mathbf{z}\,\mathbf{z}'}\left\{\exists f = (D', w') \in \mathcal{D}_{\varepsilon} \times \mathcal{W}_{\varepsilon}, \ (P_{\mathbf{z}'} - P_{\mathbf{z}})l(\cdot, f) > \varpi\right\}$$

$$\leq \left(\frac{8(r/2)^{1/(d+1)}}{\varepsilon}\right)^{(d+1)k} \exp(-m\varpi^2/(2b^2)).$$

The bound on $\Pr(J \cap \overline{Z})$ now follows.

The stage is now set to prove Theorem 3; the full proof is in Appendix D.

Proof sketch (of Theorem 3): Proposition 1 and Lemmas 5 and 6 imply that

$$\Pr_{\mathbf{z}} \left\{ \exists f \in \mathcal{F}_{\mu}, \left[\operatorname{margin}_{s}(D, \mathbf{x}) > \iota \right] \text{ and } ((P - P_{\mathbf{z}})l(\cdot, f) > t) \right\}$$

$$\leq 2 \left(\left(\frac{8(r/2)^{1/(d+1)}}{\varepsilon} \right)^{(d+1)k} \exp(-m\varpi^{2}/(2b^{2})) + \delta \right).$$

Let $s \in [k]$ and $\mu > 0$ be fixed a priori. Setting $\varepsilon = \frac{1}{m}$ in the above, elementary manipulations can show that provided $m > \frac{387}{\mathrm{margin}_s(D,\mathbf{x})^2\lambda}$, with probability at least $1 - \delta$ over $\mathbf{z} \sim P^m$, for any $f = (D, w) \in \mathcal{F}$ satisfying $\mu_s(D) \geq \mu$ and $[\mathrm{margin}_s(D, \mathbf{x}) > \iota]$, the generalization error $(P - P_{\mathbf{z}})l(\cdot, f)$ is bounded by:

$$\begin{split} &2b\sqrt{\frac{2((d+1)k\log(8m)+k\log\frac{r}{2}+\log\frac{4}{\delta})}{m}}\\ &+\frac{4L}{m}\left(\frac{1}{\lambda}\left(1+\frac{3r\sqrt{s}}{\mu}\right)\right)+\frac{2b}{m}\left(dk\log\frac{3096}{\mathrm{margin}_s^2(D,\mathbf{x})\cdot\lambda}+\log(2m+1)+\log\frac{4}{\delta}\right). \end{split}$$

It remains to distribute a prior across the bounds for each choice of s and μ . To each choice of $s \in [k]$ assign prior probability $\frac{1}{k}$. To each choice of $i \in \mathbb{N} \cup \{0\}$ for $2^{-i} \leq \mu$ assign prior probability $(i+1)^{-2}$. For a given choice of $s \in [k]$ and $2^{-i} \leq \mu$ we use $\delta(s,i) := \frac{6}{\pi^2} \frac{1}{(i+1)^2} \frac{1}{k} \delta$ (since $\sum_{i=1}^{\infty} \frac{1}{i^2} = \frac{\pi^2}{6}$). The theorem now follows.

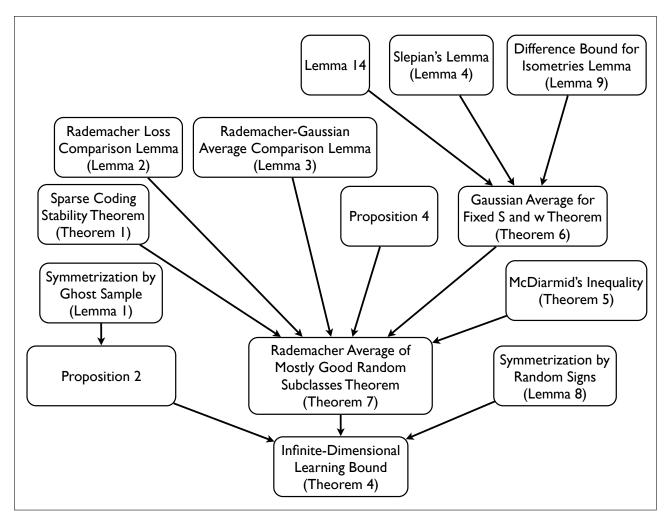


Figure 3: Proof flowchart for the Infinite-Dimensional Learning Bound (Theorem 4).

5 Infinite-dimensional setting

In the infinite-dimensional setting learning bounds with dependence on d are useless. Unfortunately, the strategy of the previous section breaks down in the infinite-dimensional setting because the straightforward construction of any epsilon-cover over the space of dictionaries had cardinality that depends on d. Even worse, epsilon-covers actually were used both to approximate the function class \mathcal{F} in $\|\cdot\|_{\infty}$ norm and to guarantee that most points of the ghost sample are good provided that all points of the training sample were good (the Good Ghost Lemma (Lemma 5)).

These issues can be overcome by requiring an additional, unlabeled sample — a device often justified in supervised learning problems because unlabeled data may be inexpensive and yet quite helpful — and by switching to more sophisticated techniques based on conditional Rademacher and Gaussian averages. After learning a hypothesis $\hat{f}_{\mathbf{z}}$ from a predictive sparse coding algorithm, the sparsity level and coding margin are measured on a second, unlabeled sample \mathbf{x}'' of m points⁷. Since this sample is independent of the choice of $\hat{f}_{\mathbf{z}}$, it is possible to guarantee that all but a very

 $^{^7\}mathrm{The}$ cardinality matches the size of the training sample $\mathbf z$ purely for simplicity.

small fraction $(\frac{\eta}{m} = \frac{\log \frac{1}{\delta}}{m})$ of points of a ghost sample **z** are good with probability $1 - \delta$. In the likely case of this good event, and for a fixed sample, we then consider all possible choices of a set of η bad indices in the ghost sample; each of the $\binom{m}{\eta}$ cases corresponds to a subclass of functions. We then approximate each subclass by a special ε -cover that is a disjoint union of a finite number of special subclasses; for each of these smaller subclasses, we bound the conditional Rademacher average by exploiting a sparsity property. The proof flowchart in Figure 3 shows the structure of the proof of Theorem 4.

5.1 Symmetrization and decomposition

The proof of the infinite-dimensional setting learning bound Theorem 4 depends critically on Lemma 9, a lemma which is non-trivial only for dictionaries with non-zero 2s-incoherence. The s-incoherence also will continue to play an important role, as it did in the overcomplete setting. Therefore, rather than wielding the deterministic subclass \mathcal{F}_{μ} of the previous section, we will work with a deterministic subclass with lower bounded s-incoherence and lower bounded 2s-incoherence.

Let $\mu^* = (\mu_s^*, \mu_{2s}^*) \in \mathbb{R}_+^2$ and define the deterministic subclass

$$\mathcal{F}_{\mu^*} = \{ f = (D, w) \in \mathcal{F} : (\mu_s(D) \ge \mu_s^*) \text{ and } (\mu_{2s}(D) \ge \mu_{2s}^*) \}.$$

The next result is immediate from Lemma 1, taking the random subclass $\mathcal{F}(\cdot)$ to be

$$\mathcal{F}(\mathbf{z}, \mathbf{x}'') := \left\{ \{\hat{f}_{\mathbf{z}}\} \cap \{f \in \mathcal{F}_{\boldsymbol{\mu}^*} : \left[\operatorname{margin}_s(D, \mathbf{x} \cup \mathbf{x}'') > \tau \right] \right\}.$$

Proposition 2 If $m \geq \left(\frac{b}{t}\right)^2$, then

$$\Pr_{\mathbf{z}\,\mathbf{x}''} \left\{ \hat{f}_{\mathbf{z}} \in \mathcal{F}_{\boldsymbol{\mu}^*} \left[\operatorname{margin}_{s}(\hat{D}_{\mathbf{z}}, \mathbf{x} \cup \mathbf{x}'') > \tau \right] \text{ and } \left((P - P_{\mathbf{z}}) l(\cdot, \hat{f}_{\mathbf{z}}) \ge t \right) \right\} \\
\leq 2 \Pr_{\mathbf{z}\,\mathbf{z}'\mathbf{x}''} \left\{ \hat{f}_{\mathbf{z}} \in \mathcal{F}_{\boldsymbol{\mu}^*} \left[\operatorname{margin}_{s}(\hat{D}_{\mathbf{z}}, (\mathbf{x} \cup \mathbf{x}'')) > \tau \right] \text{ and } \left((P_{\mathbf{z}'} - P_{\mathbf{z}}) l(\cdot, \hat{f}_{\mathbf{z}}) \ge \frac{t}{2} \right) \right\}. \quad \blacksquare$$

Now, observe that the probability of interest can be split into the probability of a large deviation happening under a "good" event and the probability of a "bad" event occurring:

$$\begin{split} & \operatorname{Pr}_{\mathbf{z}\,\mathbf{z}'\mathbf{x}''}\left\{\begin{array}{l} \hat{f}_{\mathbf{z}} \in \mathcal{F}_{\boldsymbol{\mu}^*} \left[\operatorname{margin}_s(\hat{D}_{\mathbf{z}},\mathbf{x} \cup \mathbf{x}'') > \tau\right] \text{ and } \left((P_{\mathbf{z}'} - P_{\mathbf{z}})l(\cdot,\hat{f}_{\mathbf{z}}) \geq \frac{t}{2}\right) \right. \right\} \\ & = \operatorname{Pr}_{\mathbf{z}\,\mathbf{z}'\mathbf{x}''}\left\{\begin{array}{l} \hat{f}_{\mathbf{z}} \in \mathcal{F}_{\boldsymbol{\mu}^*} \left[\operatorname{margin}_s(\hat{D}_{\mathbf{z}},\mathbf{x} \cup \mathbf{x}'') > \tau\right] \text{ and } \left(\exists \tilde{\mathbf{x}} \subseteq_{\boldsymbol{\eta}} \mathbf{x}' \left[\operatorname{margin}_s(\hat{D}_{\mathbf{z}},\tilde{\mathbf{x}}) > \tau\right]\right) \right. \\ & \left. + \operatorname{Pr}_{\mathbf{z}\,\mathbf{z}'\mathbf{x}''} \left\{\begin{array}{l} \hat{f}_{\mathbf{z}} \in \mathcal{F}_{\boldsymbol{\mu}^*} \left[\operatorname{margin}_s(\hat{D}_{\mathbf{z}},\mathbf{x} \cup \mathbf{x}'') > \tau\right] \text{ and } \left(\nexists \tilde{\mathbf{x}} \subseteq_{\boldsymbol{\eta}} \mathbf{x}' \left[\operatorname{margin}_s(\hat{D}_{\mathbf{z}},\tilde{\mathbf{x}}) > \tau\right]\right) \right. \\ & \left. + \operatorname{Pr}_{\mathbf{z}\,\mathbf{z}'\mathbf{x}''} \left\{\begin{array}{l} \exists f \in \mathcal{F}_{\boldsymbol{\mu}^*}, \left[\operatorname{margin}_s(\hat{D}_{\mathbf{z}},\mathbf{x}) > \tau\right] \text{ and } \left(\exists \tilde{\mathbf{x}} \subseteq_{\boldsymbol{\eta}} \mathbf{x}' \left[\operatorname{margin}_s(\hat{D}_{\mathbf{z}},\tilde{\mathbf{x}}) > \tau\right]\right) \right. \right. \\ & \leq \operatorname{Pr}_{\mathbf{z}\,\mathbf{z}'} \left\{\begin{array}{l} \exists f \in \mathcal{F}_{\boldsymbol{\mu}^*}, \left[\operatorname{margin}_s(\hat{D}_{\mathbf{z}},\mathbf{x}) > \tau\right] \text{ and } \left(\exists \tilde{\mathbf{x}} \subseteq_{\boldsymbol{\eta}} \mathbf{x}' \left[\operatorname{margin}_s(\hat{D}_{\mathbf{z}},\tilde{\mathbf{x}}) > \tau\right]\right) \right. \\ & \left. + \operatorname{Pr}_{\mathbf{x}'\mathbf{x}''} \left\{\begin{array}{l} \hat{f}_{\mathbf{z}} \in \mathcal{F}_{\boldsymbol{\mu}^*} \left[\operatorname{margin}_s(\hat{D}_{\mathbf{z}},\mathbf{x}'') > \tau\right] \text{ and } \left(\nexists \tilde{\mathbf{x}} \subseteq_{\boldsymbol{\eta}} \mathbf{x}' \left[\operatorname{margin}_s(\hat{D}_{\mathbf{z}},\tilde{\mathbf{x}}) > \tau\right]\right) \right. \right\}. \end{aligned} \right. \\ & \left. + \operatorname{Pr}_{\mathbf{x}'\mathbf{x}''} \left\{\begin{array}{l} \hat{f}_{\mathbf{z}} \in \mathcal{F}_{\boldsymbol{\mu}^*} \left[\operatorname{margin}_s(\hat{D}_{\mathbf{z}},\mathbf{x}'') > \tau\right] \text{ and } \left(\nexists \tilde{\mathbf{x}} \subseteq_{\boldsymbol{\eta}} \mathbf{x}' \left[\operatorname{margin}_s(\hat{D}_{\mathbf{z}},\tilde{\mathbf{x}}) > \tau\right]\right) \right. \right\}. \end{array}$$

Of the two probabilities summed in the last line, we treat the first in the next subsection. To bound the second one, note that for each choice of \mathbf{x} , $\hat{f}_{\mathbf{z}}$ is a fixed function. Hence, it is sufficient to select η such that, for any fixed function $f = (D, w) \in \mathcal{F}$, this second probability is bounded by δ . The next lemma accomplishes this bound:

Lemma 7 (Unlikely Bad Ghost) Let $f = (D, w) \in \mathcal{F}$ be fixed. If $\eta = \log \frac{1}{\delta}$, then

$$\mathrm{Pr}_{\mathbf{x}'\mathbf{x}''}\left\{ \ \left[\mathrm{margin}_s(D,\mathbf{x}'') > \tau\right] \ \mathbf{and} \ \left(\nexists \ \tilde{\mathbf{x}} \subseteq_{\eta} \mathbf{x}' \ \left[\mathrm{margin}_s(D,\tilde{\mathbf{x}}) > \tau\right] \right) \ \right\} \leq \delta.$$

Proof sketch: The proof just uses the same standard permutation argument as in the proof of the Good Ghost Lemma (Lemma 5).

5.2 Rademacher bound in the case of the good event

We now bound the probability of a large deviation in the (likely) case of the good event. Denote by $\mathcal{F}_{\mu^*}(\mathbf{x})$ the intersection of the deterministic subclass \mathcal{F}_{μ^*} with the random subclass of functions for which the Sparse Coding Stability Theorem (Theorem 1) kicks in with constants (μ_s^*, s, τ) :

$$\mathcal{F}_{\mu^*}(\mathbf{x}) := \{ f \in \mathcal{F}_{\mu^*} : [\text{margin}_s(D, \mathbf{x}) > \tau] \}.$$

This is the "good" random subclass. Similarly, let $\mathcal{F}_{\mu^*,\eta}(\mathbf{x})$ denote the "mostly good" (or "all-but- η -good") random subclass:

$$\mathcal{F}_{\boldsymbol{\mu}^*,\eta}(\mathbf{x}) := \Big\{ f \in \mathcal{F}_{\boldsymbol{\mu}^*} : \exists \, \tilde{\mathbf{x}} \subseteq_{\eta} \mathbf{x} \, \left[\operatorname{margin}_s(D, \tilde{\mathbf{x}}) > \tau \right] \Big\}.$$

Recall that $\sigma_1, \ldots, \sigma_m$ are independent Rademacher random variables.

Lemma 8 (Symmetrization by Random Signs)

$$\begin{split} & \operatorname{Pr}_{\mathbf{z}\,\mathbf{z}'}\left\{\begin{array}{l} \exists f \in \mathcal{F}_{\boldsymbol{\mu^*}}, \ \left[\operatorname{margin}_s(D,\mathbf{x}) > \tau\right] \ \mathbf{and} \ \left(\exists \ \tilde{\mathbf{x}} \subseteq_{\eta} \mathbf{x}' \ \left[\operatorname{margin}_s(D,\tilde{\mathbf{x}}) > \tau\right]\right) \\ & \mathbf{and} \ \left((P_{\mathbf{z}'} - P_{\mathbf{z}})l(\cdot,f) \geq \frac{t}{2}\right) \end{array}\right. \\ & \leq \operatorname{Pr}_{\mathbf{z},\boldsymbol{\sigma}}\left\{\sup_{f \in \mathcal{F}_{\boldsymbol{\mu^*}}(\mathbf{x})} \frac{1}{m} \sum_{i=1}^m \sigma_i l(y_i,f(x_i)) \geq \frac{t}{4}\right\} + \operatorname{Pr}_{\mathbf{z},\boldsymbol{\sigma}}\left\{\sup_{f \in \mathcal{F}_{\boldsymbol{\mu^*},\eta}(\mathbf{x})} \frac{1}{m} \sum_{i=1}^m \sigma_i l(y_i,f(x_i)) \geq \frac{t}{4}\right\}. \end{split}$$

Proof: From the definitions of the random subclasses $\mathcal{F}_{\mu^*}(\cdot)$ and $\mathcal{F}_{\mu^*,\eta}(\cdot)$, the left hand side in the lemma is equal to

$$\Pr_{\mathbf{z}\mathbf{z}'} \left\{ \sup_{f \in \mathcal{F}_{\boldsymbol{\mu}^*}(\mathbf{x}) \cap \mathcal{F}_{\boldsymbol{\mu}^*, \eta}(\mathbf{x}')} \frac{1}{m} \sum_{i=1}^m \left(l(y_i', f(x_i')) - l(y_i, f(x_i)) \right) \ge \frac{t}{2} \right\}.$$

Now, by a routine application of symmetrization by random signs this is equal to

$$\Pr_{\mathbf{z}\mathbf{z}',\sigma} \left\{ \sup_{f \in \mathcal{F}_{\boldsymbol{\mu^*}}(\mathbf{x}) \cap \mathcal{F}_{\boldsymbol{\mu^*},\eta}(\mathbf{x}')} \frac{1}{m} \sum_{i=1}^{m} \sigma_i(l(y_i', f(x_i')) - l(y_i, f(x_i))) \ge \frac{t}{2} \right\}$$

$$\leq \Pr_{\mathbf{z},\sigma} \left\{ \sup_{f \in \mathcal{F}_{\boldsymbol{\mu^*}}(\mathbf{x})} \frac{1}{m} \sum_{i=1}^{m} \sigma_i l(y_i, f(x_i)) \ge \frac{t}{4} \right\} + \Pr_{\mathbf{z},\sigma} \left\{ \sup_{f \in \mathcal{F}_{\boldsymbol{\mu^*},\eta}(\mathbf{x})} \frac{1}{m} \sum_{i=1}^{m} \sigma_i l(y_i, f(x_i)) \ge \frac{t}{4} \right\}. \quad \blacksquare$$

Since the good random subclass $\mathcal{F}_{\mu^*}(\mathbf{x})$ is just the all-but-0-good random subclass $\mathcal{F}_{\mu^*,0}(\mathbf{x})$, it is sufficient to bound the second term of the last line above for arbitrary $\eta \in [m]$. For fixed \mathbf{z} , the randomness of the subclass is annihilated and the above supremum over $\mathcal{F}_{\mu^*,\eta}(\mathbf{x})$ is a conditional Rademacher average. Bounding this conditional Rademacher average will call for a few results on the *Gaussian* average of a related function class.

First, note that for any $D \in \mathcal{D}$, the dictionary D can be factorized as D = US, where all $U \in \mathcal{U} \subset \mathbb{R}^{d \times k}$ satisfy the isometry property $U^TU = I$, and S lives in a space $S := (B_{\mathbb{R}^k})^k$ of lower-dimensional dictionaries (Maurer and Pontil, 2010). Consider a particular choice of $S \in S$, linear hypothesis $w \in \mathcal{W}$, and m-sample \mathbf{x} . The subclass of interest will be those functions corresponding to $U \in \mathcal{U}$ such that the encoder φ_{US} is s-sparse on \mathbf{x} . It turns out that the Gaussian average of this subclass is well-behaved.

Recall that $\gamma = (\gamma_1, \dots \gamma_m)$ where the γ_i are iid standard normals.

Theorem 6 (Gaussian Average for Fixed S and w) Let $S \in S$, $s \in [k]$, and x be a fixed m-sample. Denote by \mathcal{U}_x the particular subclass of \mathcal{U} defined as:

$$\mathcal{U}_{\mathbf{x}} := \{ U \in \mathcal{U} : s\text{-sparse}(\varphi_{US}(\mathbf{x})) \}.$$

Then

$$\mathsf{E}_{\gamma} \sup_{U \in \mathcal{U}_{\mathbf{x}}} \frac{2}{m} \sum_{i=1}^{m} \gamma_i \langle w, \varphi_{US}(x_i) \rangle \le \frac{4rk\sqrt{2s}}{\mu_{2s}(S)\sqrt{m}}. \tag{15}$$

The proof of this result uses the following lemma that shows how the difference between the feature maps φ_{US} and $\varphi_{U'S}$ can be characterized by the difference between U and U'. Define the s-restricted 2-norm of S as $||S||_{2,s} := \sup_{t \in \mathbb{R}^n: ||t||=1, |\sup_{t \in \mathbb{R}^n: ||t||=1, |\sup_{t \in \mathbb{R}^n: |t||=1, |t|=1, |t$

Lemma 9 (Difference Bound for Isometries) Let $U, U' \in \mathcal{U}$ be isometries as above, $S \in \mathcal{S}$, and $x \in B_{\mathbb{R}^d}$. If $\|\varphi_{US}(x)\|_0 \leq s$ and $\|\varphi_{U'S}(x)\|_0 \leq s$, then

$$\|\varphi_{US}(x) - \varphi_{U'S}(x)\|_{2} \le \frac{2\|S\|_{2,2s}}{\mu_{2s}(S)} \|(U'^{T} - U^{T})x\|_{2}.$$

Proof sketch: The proof uses a perturbation analysis of solutions to linearly constrained positive definite quadratic programs (Daniel, 1973), exploiting the sparsity of the optimal solutions to have dependence only on $||S||_{2,2s}$ and $\mu_{2s}(S)$ rather than $||S||_2$ and $\mu_k(S)$.

Proof (of Theorem 6): Define a Gaussian process Ω , indexed by U, by $\Omega_U := \sum_{i=1}^m \gamma_i \langle w, \varphi_{US}(x_i) \rangle$. Our goal is to apply Slepian's Lemma (Lemma 4) to bound the expectation of the supremum of Ω , which depends on φ_{US} , by the expectation of the supremum of a Gaussian process Γ which depends only on U.

$$\mathsf{E}_{\gamma} \left(\Omega_{U} - \Omega_{U'}\right)^{2} = \mathsf{E}_{\gamma} \left(\sum_{i=1}^{m} \gamma_{i} \left\langle w, \varphi_{US}(x_{i}) \right\rangle - \sum_{i=1}^{m} \gamma_{i} \left\langle w, \varphi_{U'S}(x_{i}) \right\rangle \right)^{2}$$

$$= \sum_{i=1}^{m} \left(\left\langle w, \varphi_{US}(x_{i}) - \varphi_{U'S}(x_{i}) \right\rangle \right)^{2}$$

$$\leq r^{2} \sum_{i=1}^{m} \|\varphi_{US}(x_{i}) - \varphi_{U'S}(x_{i})\|^{2}$$

$$(16)$$

Applying the result from Lemma 9, we have

$$\begin{split} \mathsf{E}_{\gamma} \left(\Omega_{U} - \Omega_{U'} \right)^{2} &\leq r^{2} \sum_{i=1}^{m} \| \varphi_{US}(x_{i}) - \varphi_{U'S}(x_{i}) \|^{2} \\ &\leq \left(\frac{2r \| S \|_{2,2s}}{\mu_{2s}(S)} \right)^{2} \sum_{i=1}^{m} \left\| \left(U'^{T} - U^{T} \right) x_{i} \right\|_{2}^{2} \\ &= \left(\frac{2r \| S \|_{2,2s}}{\mu_{2s}(S)} \right)^{2} \sum_{i=1}^{m} \sum_{j=1}^{k} \left(\langle U'e_{j}, x_{i} \rangle - \langle Ue_{j}, x_{i} \rangle \right)^{2} \\ &= \left(\frac{2r \| S \|_{2,2s}}{\mu_{2s}(S)} \right)^{2} \mathsf{E}_{\gamma} \left(\left(\sum_{i=1}^{m} \sum_{j=1}^{k} \gamma_{ij} \langle U'e_{j}, x_{i} \rangle \right) - \left(\sum_{i=1}^{m} \sum_{j=1}^{k} \gamma_{ij} \langle Ue_{j}, x_{i} \rangle \right) \right)^{2} \\ &= \mathsf{E}_{\gamma} \left(\Gamma_{U} - \Gamma_{U'} \right)^{2} \end{split}$$

for

$$\Gamma_U := \frac{2r||S||_{2,2s}}{\mu_{2s}(S)} \sum_{i=1}^m \sum_{j=1}^k \gamma_{ij} \langle Ue_j, x_i \rangle.$$

By Slepian's Lemma (Lemma 4), $\mathsf{E}_{\gamma}\sup_{U}\Omega_{U} \leq \mathsf{E}_{\gamma}\sup_{U}\Gamma_{U}$. It remains to bound $E_{\gamma}\sup_{U}\Gamma_{U}$:

$$\frac{\mu_{2s}(S)}{2r\|S\|_{2,2s}} \mathsf{E}_{\gamma} \sup_{U} \Gamma_{U} = \mathsf{E}_{\gamma} \sup_{U} \sum_{i=1}^{m} \sum_{j=1}^{k} \gamma_{ij} \langle Ue_{j}, x_{i} \rangle$$

$$= \mathsf{E}_{\gamma} \sup_{U} \sum_{j=1}^{k} \langle Ue_{j}, \sum_{i=1}^{m} \gamma_{ij} x_{i} \rangle$$

$$\leq \mathsf{E}_{\gamma} \sup_{U} \sum_{j=1}^{k} \|Ue_{j}\| \|\sum_{i=1}^{m} \gamma_{ij} x_{i}\|$$

$$= k \mathsf{E}_{\gamma} \|\sum_{i=1}^{m} \gamma_{i1} x_{i}\|$$

$$\leq k \sqrt{\mathsf{E}_{\gamma} \|\sum_{i=1}^{m} \gamma_{i1} x_{i}\|^{2}}$$

$$= k \sqrt{\mathsf{E}_{\gamma} \left\langle \sum_{i=1}^{m} \gamma_{i1} x_{i}, \sum_{i=1}^{m} \gamma_{i1} x_{i} \right\rangle} = k \sqrt{\sum_{i=1}^{m} \|x_{i}\|^{2}} \leq k \sqrt{m}.$$

Hence,

$$\mathsf{E}_{\gamma} \sup_{U \in \mathcal{U}} \frac{2}{m} \sum_{i=1}^m \gamma_i \langle w, \varphi_{US}(x_i) \rangle \quad \leq \quad \frac{4r \|S\|_{2,2s} k}{\mu_{2s}(S) \sqrt{m}} \quad \leq \quad \frac{4r k \sqrt{2s}}{\mu_{2s}(S) \sqrt{m}},$$

where we used the fact that $||S||_{2,2s} \leq \sqrt{2s}$ (see Lemma 14 in Appendix E for a proof).

We present the *pièce de résistance* of this section:

Theorem 7 (Rademacher Average of Mostly Good Random Subclasses)

$$\Pr_{\mathbf{z}, \boldsymbol{\sigma}} \left\{ \sup_{f \in \mathcal{F}_{\boldsymbol{\mu}^*, \boldsymbol{\eta}}(\mathbf{x})} \frac{1}{m} \sum_{i=1}^m \sigma_i l(y_i, f(x_i)) \ge \frac{t}{4} \right\} \le {m \choose \eta} \left(\frac{8(r/2)^{1/(k+1)}}{\varepsilon} \right)^{(k+1)k} \exp(-mt_3^2/(2b^2)),$$

for

$$t_3 := \frac{t}{4} - \frac{L\varepsilon}{\lambda} \left(\frac{3r\sqrt{s}}{\mu_s^*} + 1 \right) - \frac{2L\sqrt{\pi}rk\sqrt{s}}{\mu_{2s}^*\sqrt{m}} - \frac{2b\eta}{m}.$$

Proof: As before, each dictionary $D \in \mathcal{D}$ will be factorized as D = US for U an isometry in \mathcal{U} and $S \in \mathcal{S} = (B_{\mathbb{R}^k})^k$. Let $\mathcal{S}_{\varepsilon}$ be a minimum-cardinality proper ε -cover (in operator norm) of $\{S \in \mathcal{S} : \mu_s(S) \geq \mu_s^*, \mu_{2s}(S) \geq \mu_{2s}^*\}$, the set of suitably incoherent elements of \mathcal{S} .

Recall that the goal is to control the Rademacher complexity of $\mathcal{F}_{\mu^*,\eta}(\mathbf{x})$. Our strategy will be to control this complexity by controlling the complexity of each subclass from a partition of $\mathcal{F}_{\mu^*,\eta}(\mathbf{x})$. For an arbitrary $f = (D, w) \in \mathcal{F}$, let an index i be good if and only if $[\text{margin}_s(D, x_i) > \tau]$, and let an index be bad if and only if it is not good. Consider a fixed m-sample \mathbf{z} and the occurrence of a set of $m - \eta$ good indices⁸. There are $N := \binom{m}{\eta}$ ways to choose this set of indices. We can partition $\mathcal{F}_{\mu^*,\eta}(\mathbf{x})$ into N subclasses $\mathcal{F}^1_{\mu^*,\eta}(\mathbf{x}),\ldots,\mathcal{F}^N_{\mu^*,\eta}(\mathbf{x})$ such that for all functions in a given subclass, a particular set of $m - \eta$ indices is guaranteed to be good. To be precise, we can choose distinct good index sets Γ_1,\ldots,Γ_N , each of cardinality $m - \eta$, such that for each Γ_j , if $i \in \Gamma_j$ then all f = (D, w) in $\mathcal{F}^j_{\mu^*,\eta}$ satisfy $[\text{margin}_s(D, x_i) > \tau]$.

Since the $\mathcal{F}_{\mu^*,\eta}^j(\mathbf{x})$ form a partition, we can control the complexity of $\mathcal{F}_{\mu^*,\eta}(\mathbf{x})$ via:

$$\sup_{f \in \mathcal{F}_{\boldsymbol{\mu^*}, \eta}(\mathbf{x})} \sum_{i=1}^m \sigma_i l(y_i, f(x_i)) = \max_{j \in [N]} \sup_{f \in \mathcal{F}_{\boldsymbol{\mu^*}, \eta}^j(\mathbf{x})} \sum_{i=1}^m \sigma_i l(y_i, f(x_i)).$$

To gain a handle on the complexity of each subclass $\mathcal{F}^{j}_{\mu^*,\eta}(\mathbf{x})$, we will approximate the subclasses as follows. For each $j \in [N]$, define an ε -neighborhood of $\mathcal{F}^{j}_{\mu^*,\eta}(\mathbf{x})$ as

$$\bar{\mathcal{F}}_{\mu^*,\eta}^j(\mathbf{x}) := \left\{ \begin{array}{ccc} f = (US',w'): & \|S-S'\| \leq \varepsilon, & \|w-w'\| \leq \varepsilon, \\ & S \in \mathcal{S}, & w \in \mathcal{W}, \end{array} \right. \quad (US,w) \in \mathcal{F}_{\boldsymbol{\mu^*},\eta}^j(\mathbf{x}) \ \left. \right\};$$

note that the ε neighborhood is taken with respect to S and w but not U. Also, let $\mathcal{W}_{\varepsilon}$ be a minimum-cardinality ε -cover of \mathcal{W} and define an infinite-cardinality epsilon-cover of \mathcal{F} :

$$\mathcal{F}_{\varepsilon} := \left\{ f = (US', w') \in \mathcal{F} : U \in \mathcal{U}, S' \in \mathcal{S}_{\varepsilon}, w' \in \mathcal{W}_{\varepsilon} \right\}.$$

Finally, taking the intersection of $\bar{\mathcal{F}}^{j}_{\mu^*,\eta}(\mathbf{x})$ with $\mathcal{F}_{\varepsilon}$ yields the $\mathcal{F}^{j}_{\mu^*,\eta}(\mathbf{x})$ -approximating subclass, a disjoint union of subclasses equal to

$$\bigcup_{S' \in \mathcal{S}_{\varepsilon}, w' \in \mathcal{W}_{\varepsilon}} \mathcal{F}_{\boldsymbol{\mu^*}, \eta}^{j, S', w'}(\mathbf{x})$$

⁸Each of the remaining indices can be either good or bad.

for

$$\mathcal{F}_{\boldsymbol{\mu}^*,\eta}^{j,S',w'}(\mathbf{x}) := \mathcal{F}_{\boldsymbol{\mu}^*,\eta}^{j}(\mathbf{x}) \cap \left\{ f \in \mathcal{F} : f = (US',w') : U \in \mathcal{U} \right\}.$$

To show that this disjoint union is a good approximator for $\mathcal{F}_{\mu^*,\eta}^j(\mathbf{x})$, for each $j \in [N]$ and arbitrary $\boldsymbol{\sigma} \in \{-1,1\}^m$ we compare

$$\sup_{f \in \mathcal{F}_{\boldsymbol{\mu}^*, \eta}^{j}(\mathbf{x})} \frac{1}{m} \sum_{i=1}^{m} \sigma_{i} l(y_{i}, f(x_{i})) \quad \text{and} \quad \max_{S' \in \mathcal{S}_{\varepsilon}, w' \in \mathcal{W}_{\varepsilon}} \sup_{f \in \mathcal{F}_{\boldsymbol{\mu}^{*}, \eta}^{j, S', w'}(\mathbf{x})} \sum_{i=1}^{m} \sigma_{i} l(y_{i}, f(x_{i})).$$

Without loss of generality, choose j=1 and take $\Gamma_1=[m-\eta]$. If f is in $\mathcal{F}^1_{\mu^*,\eta}(\mathbf{x})$, it follows that there exists an f' in the disjoint union $\bigcup_{S'\in\mathcal{S}_{\varepsilon},w'\in\mathcal{W}_{\varepsilon}}\mathcal{F}^{1,S',w'}_{\mu^*,\eta}(\mathbf{x})$ such that

$$\frac{1}{m} \sum_{i=1}^{m} \sigma_{i} \left| l(y_{i}, \langle w, \varphi_{D}(x_{i}) \rangle) - l(y_{i}, \langle w', \varphi_{D'}(x_{i}) \rangle) \right| \\
\leq \frac{L}{m} \left(\sum_{i=1}^{m-\eta} \sigma_{i} \left| \langle w, \varphi_{D}(x_{i}) \rangle - \langle w', \varphi_{D'}(x_{i}) \rangle \right| \right) + \frac{1}{m} \sum_{i=m-\eta+1}^{m} \sigma_{i} \left| l(y_{i}, \langle w, \varphi_{D}(x_{i}) \rangle) - l(y_{i}, \langle w', \varphi_{D'}(x_{i}) \rangle) \right| \\
\leq \frac{L\varepsilon}{\lambda} \left(\frac{3r\sqrt{s}}{\mu_{s}^{*}} + 1 \right) + \frac{b\eta}{m},$$

where the last line is due to the Sparse Coding Stability Theorem (Theorem 1).

Therefore, for any $\sigma \in \{-1,1\}^m$ it holds that

$$\sup_{f \in \mathcal{F}_{\mu^*,\eta}(\mathbf{x})} \frac{1}{m} \sum_{i=1}^{m} \sigma_i l(y_i, f(x_i))$$

$$\leq \max_{j \in [N]} \max_{S' \in \mathcal{S}_{\varepsilon}, w' \in \mathcal{W}_{\varepsilon}} \sup_{f \in \mathcal{F}_{\mu^*,\eta}^{j,S',w'}(\mathbf{x})} \frac{1}{m} \sum_{i=1}^{m} \sigma_i l(y_i, f(x_i)) + \frac{L\varepsilon}{\lambda} \left(\frac{3r\sqrt{s}}{\mu_s^*} + 1 \right) + \frac{b\eta}{m}.$$

Thus, the approximation error from using the disjoint union is small (it is $O(\frac{1}{m})$ if $\varepsilon = \frac{1}{m}$).

It remains to control the complexity of the approximating subclass. From the above, for fixed **z**:

$$\Pr_{\sigma} \left\{ \sup_{f \in \mathcal{F}_{\mu^*,\eta}(\mathbf{x})} \frac{1}{m} \sum_{i=1}^{m} \sigma_i l(y_i, f(x_i)) \ge \frac{t}{4} \right\}$$

$$\leq \Pr_{\sigma} \left\{ \max_{\substack{j \in [N] \\ S' \in \mathcal{S}_{\varepsilon}, w' \in \mathcal{W}_{\varepsilon}}} \sup_{f \in \mathcal{F}_{\mu^*,\eta}^{j,S',w'}(\mathbf{x})} \frac{1}{m} \sum_{i=1}^{m} \sigma_i l(y_i, f(x_i)) \ge \frac{t}{4} - \frac{L\varepsilon}{\lambda} \left(\frac{3r\sqrt{s}}{\mu_s^*} + 1 \right) - \frac{b\eta}{m} \right\}$$

$$\leq N \left(\frac{8(r/2)^{1/(k+1)}}{\varepsilon} \right)^{(k+1)k} \cdot \frac{1}{\varepsilon} \left\{ \sup_{f \in \mathcal{F}_{\mu^*,\eta}^{j,S',w'}(\mathbf{x})} \frac{1}{m} \sum_{i=1}^{m} \sigma_i l(y_i, f(x_i)) \ge \frac{t}{4} - \frac{L\varepsilon}{\lambda} \left(\frac{3r\sqrt{s}}{\mu_s^*} + 1 \right) - \frac{b\eta}{m} \right\}.$$

Now, from McDiarmid's inequality (Theorem 5), for any fixed $j \in [N]$, $S' \in \mathcal{S}_{\varepsilon}$ and $w' \in \mathcal{W}_{\varepsilon}$,

$$\Pr_{\boldsymbol{\sigma}} \left\{ \sup_{f \in \mathcal{F}_{\boldsymbol{\mu^*}, \eta}^{j, S', w'}(\mathbf{x})} \frac{1}{m} \sum_{i=1}^{m} \sigma_i l(y_i, f(x_i)) > \mathsf{E}_{\boldsymbol{\sigma}} \sup_{f \in \mathcal{F}_{\boldsymbol{\mu^*}, \eta}^{j, S', w'}(\mathbf{x})} \frac{1}{m} \sum_{i=1}^{m} \sigma_i l(y_i, f(x_i)) + t_1 \right\}$$

is at most $\exp(-mt_1^2/(2b^2))$.

To make the above useful, let us get a handle on the Rademacher complexity term

$$\mathsf{E}_{\boldsymbol{\sigma}} \sup_{f \in \mathcal{F}_{\boldsymbol{\mu}^*, n'}^{j, S', w'}(\mathbf{x})} \frac{1}{m} \sum_{i=1}^{m} \sigma_i l(y_i, f(x_i)).$$

Without loss of generality, again take j = 1 and $\Gamma_1 = [m - \eta]$. Then

$$\mathbb{E}_{\boldsymbol{\sigma}} \sup_{f \in \mathcal{F}_{\boldsymbol{\mu}^{*}, \eta}^{j, S', w'}(\mathbf{x})} \frac{1}{m} \sum_{i=1}^{m} \sigma_{i} l(y_{i}, f(x_{i}))$$

$$\leq \mathbb{E}_{\sigma_{1}, \dots, \sigma_{m-\eta}} \left\{ \sup_{f \in \mathcal{F}_{\boldsymbol{\mu}^{*}, \eta}^{j, S', w'}(\mathbf{x})} \frac{1}{m} \sum_{i=1}^{m-\eta} \sigma_{i} l(y_{i}, f(x_{i})) \right\}$$

$$+ \mathbb{E}_{\sigma_{m-\eta+1}, \dots, \sigma_{m}} \left\{ \sup_{f \in \mathcal{F}_{\boldsymbol{\mu}^{*}, \eta}^{j, S', w'}(\mathbf{x})} \frac{1}{m} \sum_{i=m-\eta+1}^{m} \sigma_{i} l(y_{i}, f(x_{i})) \right\}$$

$$\leq \mathbb{E}_{\sigma_{1}, \dots, \sigma_{m-\eta}} \left\{ \sup_{f \in \mathcal{F}_{\boldsymbol{\mu}^{*}, \eta}^{j, S', w'}(\mathbf{x})} \frac{1}{m} \sum_{i=1}^{m-\eta} \sigma_{i} l(y_{i}, f(x_{i})) \right\} + \frac{b\eta}{m}.$$

Now, Theorem 6, the Rademacher Loss Comparison Lemma (Lemma 2), and the Rademacher-Gaussian Average Comparison Lemma (Lemma 3) imply that

$$\begin{split} \mathsf{E}_{\pmb{\sigma}} \sup_{\{U \in \mathcal{U}: s\text{-sparse}(\varphi_{US}(\mathbf{x}))\}} \frac{1}{m} \sum_{i=1}^{m-\eta} \sigma_i l\big(y_i, \langle w, \varphi_{US}(x_i) \rangle\big) &\leq \frac{\sqrt{m-\eta}}{m} \frac{2L\sqrt{\pi} r k \sqrt{s}}{\mu_{2s}(S)} \\ &\leq \frac{2L\sqrt{\pi} r k \sqrt{s}}{\mu_{2s}(S)\sqrt{m}}, \end{split}$$

and hence

$$\Pr_{\sigma} \left\{ \sup_{f \in \mathcal{F}_{\mu^*, \eta}^{j, S', w'}(\mathbf{x})} \frac{1}{m} \sum_{i=1}^{m} \sigma_i l(f(x_i)) > \frac{2L\sqrt{\pi}rk\sqrt{s}}{\mu_{2s}^* \sqrt{m}} + \frac{b\eta}{m} + t_1 \right\} \le \exp(-mt_1^2/(2b^2)).$$

Combining this bound with the fact that the bound is independent of the draw of \mathbf{z} and applying Proposition 4 (with d set to k) to extend the bound over all choices of j, S', and w' yields the final result.

For the case of $\eta = 0$, let

$$t_2 := \frac{t}{4} - \frac{L\varepsilon}{\lambda} \left(\frac{3r\sqrt{s}}{\mu_s^*} + 1 \right) - \frac{2L\sqrt{\pi}rk\sqrt{s}}{\mu_{2s}^*\sqrt{m}}.$$

It is now possible to prove the generalization error bound for the infinite-dimensional setting.

Proof (of Theorem 4): Since $\mathcal{F}_{\mu_*}(\mathbf{x})$ is equivalent to $\mathcal{F}_{\mu^*,0}(\mathbf{x})$, Lemma 8 and Theorem 7 imply that

and consequently the full probability (14) in Proposition 2 can be upper bounded (using $\eta = \log \frac{1}{\delta}$) as:

$$\begin{split} & \Pr_{\mathbf{z}\,\mathbf{x}''}\left\{ \begin{array}{l} \hat{f}_{\mathbf{z}} \in \mathcal{F}_{\boldsymbol{\mu}^*}, \ \left[\operatorname{margin}_s(\hat{D}_{\mathbf{z}}, (\mathbf{x} \cup \mathbf{x}'')) > \tau \right] \ \mathbf{and} \ \left((P - P_{\mathbf{z}}) l(\cdot, \hat{f}_{\mathbf{z}}) \geq t \right) \end{array} \right\} \\ & \leq 4 \binom{m}{\log \frac{1}{\delta}} \left(\frac{8(r/2)^{1/(k+1)}}{\varepsilon} \right)^{(k+1)k} \exp(-mt_3^2/(2b^2)) + 2\delta. \end{split}$$

After some elementary manipulations and choosing $\varepsilon = \frac{1}{m}$, we nearly have the final learning bound. Let $\mu_s^*, \mu_{2s}^* > 0$, $s \in [k]$, and $m \ge \frac{43}{\tau^2 \lambda}$ be fixed a priori. With probability at least $1 - \delta$ over a labeled m-sample $\mathbf{z} \sim P^m$ and a second, unlabeled m-sample $\mathbf{x}'' \sim \Pi^m$, if an algorithm learns hypothesis $\hat{f}_{\mathbf{z}} = (\hat{D}_{\mathbf{z}}, \hat{w}_{\mathbf{z}})$ from \mathbf{z} such that $\mu_{2s}(\hat{D}_{\mathbf{z}}) \ge \mu_{2s}^*, \mu_s(\hat{D}_{\mathbf{z}}) \ge \mu_s^*$, s-sparse($\varphi_{\hat{D}_{\mathbf{z}}}(\mathbf{x} \cup \mathbf{x}'')$), and $[\text{margin}_s(\hat{D}_{\mathbf{z}}, \mathbf{x} \cup \mathbf{x}'') > \tau]$ all hold, then the generalization error $(P - P_{\mathbf{z}})l(\cdot, \hat{f}_{\mathbf{z}})$ is bounded by:

$$\frac{8L\sqrt{\pi}rk\sqrt{s}}{\mu_{2s}^*\sqrt{m}} + b\sqrt{\frac{8\left((k+1)k\log(8m) + k\log\frac{r}{2} + (\log m + 1)\log\frac{4}{\delta} + \log 2\right)}{m}} + \frac{1}{m}\left(\frac{4L}{\lambda}\left(\frac{3r\sqrt{s}}{\mu_s^*} + 1\right) + 8b\log\frac{4}{\delta}\right).$$

Making this bound made adaptive to the incoherences, sparsity level, and margin on \mathbf{z} and \mathbf{x}'' yields the following final result. With probability at least $1 - \delta$ over a labeled m-sample $\mathbf{z} \sim P^m$ and a second, unlabeled sample $\mathbf{x}'' \sim \Pi^m$, if an algorithm learns hypothesis $\hat{f}_{\mathbf{z}} = (\hat{D}_{\mathbf{z}}, \hat{w}_{\mathbf{z}})$ such that $\varphi_{\hat{D}_{\mathbf{z}}}$ is s-sparse on $(\mathbf{x} \cup \mathbf{x}'')$, $\mu_{2s}(\hat{f}_{\mathbf{z}}) > 0$, and

$$m \ge \frac{43}{\operatorname{margin}_s^2(\hat{D}_{\mathbf{z}}, \mathbf{x} \cup \mathbf{x}'') \cdot \lambda},$$

then the generalization error $(P - P_z)l(\cdot, \hat{f}_z)$ is bounded by:

$$\frac{16L\sqrt{\pi}rk\sqrt{s}}{\mu_{2s}(\hat{f}_{\mathbf{z}})\sqrt{m}} + b\sqrt{\frac{8\left((k^2 + k)\log(8m) + k\log\frac{r}{2} + (\log m + 1)\log\frac{7\alpha k}{\delta} + \log 2\right)}{m}} + \frac{1}{m}\left(\frac{4L}{\lambda}\left(\frac{6r\sqrt{s}}{\mu_s(\hat{f}_{\mathbf{z}})} + 1\right) + 8b\log\frac{7\alpha k}{\delta}\right),$$
for $\alpha = \left(\log_2\left(\frac{4}{\mu_s(\hat{f}_{\mathbf{z}})}\right)\log_2\left(\frac{4}{\mu_s(\hat{f}_{\mathbf{z}})}\right)\right)^2$.

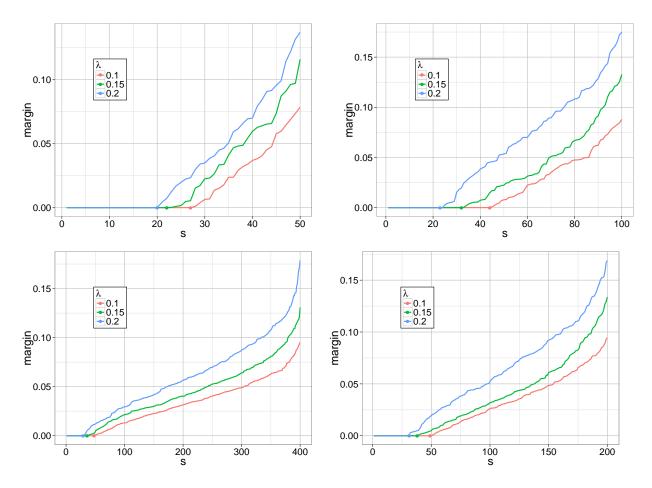


Figure 4: The s-margin for predictive sparse coding trained on the USPS training set, digit 4 versus all, for three settings of λ . Clockwise from top left: 50 atoms, 100 atoms, 200 atoms, and 400 atoms. The sparsity level (maximum number of non-zeros per code, taken across all codes of the training points) is indicated by the dots.

6 An empirical study of the s-margin

We now establish some empirical evidence that the s-margin is well away from zero even when s is only slightly larger than the observed sparsity level. We performed experiments on two separate digit classification tasks, from the USPS dataset and the MNIST dataset LeCun et al. (1998). In both cases, we employed the single binary classification task of the digit 4 versus all the other digits, and for both datasets all the training data was used. Predictive sparse coding was trained as per the stochastic gradient descent approach of Mairal et al. (2012).

The results for USPS and MNIST are shown in Figures 4 and 5 respectively. Each data point (an image) was normalized to unit norm. In all plots, it is apparent that when the minimum sparsity level is s (indicated by the colored dots on the x-axis of the plots), there is a non-trivial $(s+\rho)$ -margin for ρ a small positive integer. Using the 2s-margin when s-sparsity holds may ensure that there is a moderate margin for only a constant factor increase to s.

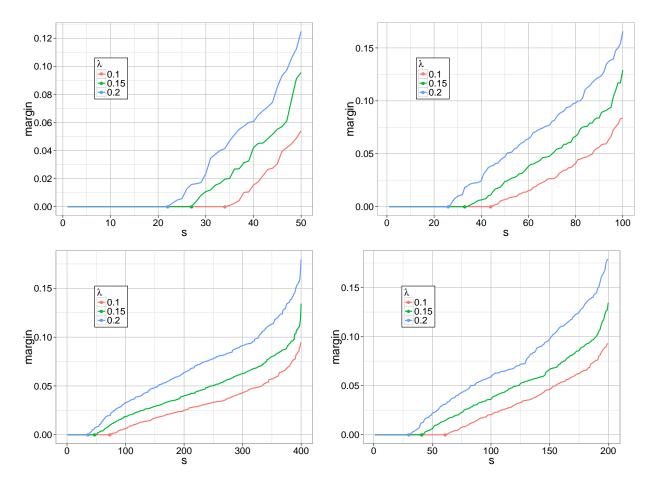


Figure 5: The s-margin for predictive sparse coding trained on the MNIST training set, digit 4 versus all, for three settings of λ . Clockwise from top left: 50 atoms, 100 atoms, 200 atoms, and 400 atoms. The sparsity level (maximum number of non-zeros per code, taken across all codes of the training points) is indicated by the dots.

7 Discussion and open problems

We have shown the first generalization error bounds for predictive sparse coding. The learning bounds in Theorems 3 and 4 are intimately related to the stability of the sparse encoder, and consequently the bounds depend on properties that depend both on the learned dictionary and the training sample. Using the techniques of this work, in the infinite-dimensional setting it is unclear whether one can achieve the encoder stability guarantees without measuring properties of the encoder on an independent, unlabeled sample. It is an important open problem whether there is a generalization error bound for the infinite-dimensional setting which does not rely on the second sample. Additionally, the PRP condition in the Sparse Coding Stability Theorem (Theorem 1) appears to be much stronger than what should be required. We conjecture that the PRP should actually be $O(\varepsilon)$ rather than $O(\sqrt{\varepsilon})$. If this conjecture turns out to be true, then the number of samples required before Theorems 3 and 4 kick in would be greatly reduced, as would be the size of many of the constants in the results.

While this work establishes upper bounds on the generalization error for predictive sparse coding, two things remain unclear. How close are these bounds to the optimal ones? Also, what lower bounds can be established in each of the settings? If the conditions on which these bounds rely are of fundamental importance, then the presented data-dependent bounds provide motivation for an algorithm to prefer dictionaries for which small subdictionaries are well-conditioned and to additionally encourage large coding margin on the training sample.

A Proof of Sparse Coding Stability Theorem

The flow of this section is as follows. We first establish some preliminary notation and summarize important conditions. Several lemmas are then presented to support a key sparsity lemma. This sparsity lemma establishes that the solution to the perturbed problem is sparse provided the perturbation is not too large. Finally, the sparsity of this new solution is exploited to bound the difference of the new solution from the old solution. This flow is embodied by the proof flowchart in Figure 6.

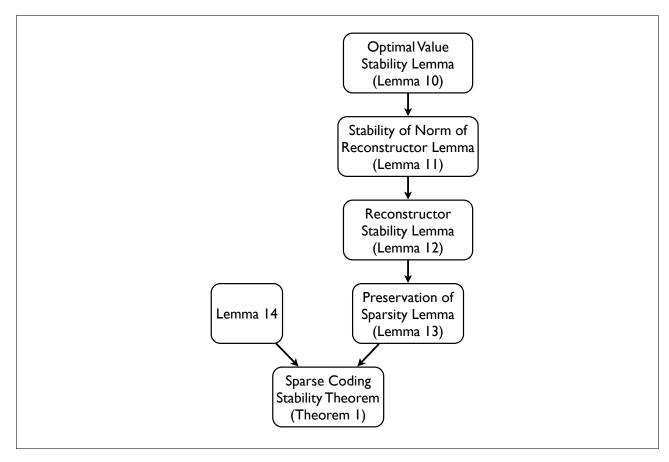


Figure 6: Proof flowchart for the Sparse Coding Stability Theorem (Theorem 1).

A.1 Notation and assumptions

Let α and $\tilde{\alpha}$ respectively denote the solutions to the LASSO problems:

$$\alpha = \arg\min_{z} \frac{1}{2} \|x - Dz\|_{2}^{2} + \lambda \|z\|_{1} \qquad \qquad \tilde{\alpha} = \arg\min_{z} \frac{1}{2} \|x - \tilde{D}z\|_{2}^{2} + \lambda \|z\|_{1}.$$

First, let's review the optimality conditions for the LASSO (Asif and Romberg, 2010, conditions

L1 and L2):

$$\langle D_j, x - D\alpha \rangle = \operatorname{sign}(\alpha_j)\lambda \quad \text{if } \alpha_j \neq 0,$$

 $|\langle D_j, x - D\alpha \rangle| < \lambda \quad \text{otherwise.}$

Note that the above optimality conditions imply that if $\alpha_j \neq 0$ then

$$|\langle D_j, x - D\alpha \rangle| = \lambda.$$

Assumptions

The statement of the Sparse Coding Stability Theorem (Theorem 1) makes the following assumptions:

(A1) - Closeness D and \tilde{D} are close, as measured by operator norm:

$$\|\tilde{D} - D\|_2 \le \varepsilon.$$

(A2) - Incoherence There is a $\mu > 0$ such that, for all $J \subseteq [k]$ satisfying |J| = s:

$$\sigma_{\min}(D_J) \geq \mu$$
.

(A3) - Sparsity with margin For some fixed $\tau > 0$, there is a $\mathcal{I} \subseteq [k]$ with $|\mathcal{I}| = k - s$ such that for all $i \in \mathcal{I}$:

$$|\langle D_i, x - D\alpha \rangle| < \lambda - \tau.$$

Consequently, all $i \in \mathcal{I}$ satisfy $\alpha_i = 0$.

A.2 Useful observations

Let v_D^* be the optimal value of the LASSO for dictionary D:

$$v_D^* = \min_{z} \frac{1}{2} ||x - Dz||_2^2 + \lambda ||z||_1$$
$$= \frac{1}{2} ||x - D\alpha||_2^2 + \lambda ||\alpha||_1$$

Likewise, let

$$v_{\tilde{D}}^* = \frac{1}{2} \|x - \tilde{D}\tilde{\alpha}\|_2^2 + \lambda \|\tilde{\alpha}\|_1$$

The first observation is that the values of the optimal solutions are close:

Lemma 10 (Optimal Value Stability) If $||D - \tilde{D}||_2 \le \varepsilon$, then

$$\left|v_D^* - v_{\tilde{D}}^*\right| \le \frac{\varepsilon}{\lambda}.$$

Proof: The proof is simple:

$$\begin{split} v_{\tilde{D}}^* &\leq \frac{1}{2} \|x - \tilde{D}\alpha\|_2^2 + \lambda \|\alpha\|_1 \\ &= \frac{1}{2} \|x - D\alpha + (D - \tilde{D})\alpha\|_2^2 + \lambda \|\alpha\|_1 \\ &\leq \frac{1}{2} \left(\|x - D\alpha\|_2^2 + \|x - D\alpha\|_2 \|(D - \tilde{D})\alpha\|_2 + \|(D - \tilde{D})\alpha\|_2^2 \right) + \lambda \|\alpha\|_1 \\ &\leq \frac{1}{2} \|x - D\alpha\|_2^2 + \lambda \|\alpha\|_1 + \frac{1}{2} \left(\frac{\varepsilon}{\lambda} + \left(\frac{\varepsilon}{\lambda} \right)^2 \right) \\ &\leq v_D^* + \frac{\varepsilon}{\lambda} \end{split}$$

for $\frac{\varepsilon}{\lambda} \leq 1$. A symmetric argument shows that $v_D^* \leq v_{\tilde{D}}^* + \frac{\varepsilon}{\lambda}$.

The second observation shows that the norms of the optimal reconstructors are close.

Lemma 11 (Stability of Norm of Reconstructor) If $||D - \tilde{D}||_2 \le \varepsilon$, then

$$\left| \|D\alpha\|_2^2 - \|\tilde{D}\tilde{\alpha}\|_2^2 \right| \le \frac{2\varepsilon}{\lambda}.$$

Showing this is more involved than the previous observation.

Proof: First, we claim (and show) that

$$(x - D\alpha)^T D\alpha = \lambda \|\alpha\|_1. \tag{17}$$

The proof of the claim comes directly from Osborne et al. (2000, circa (2.8)) To see (17), consider the equivalent dual problem for the LASSO (for appropriate choice of t):

minimize
$$\frac{1}{2} ||x - Dz||_2^2$$

subject to
$$||z||_1 \le t.$$

The Lagrangian is

$$\mathcal{L}(z,\lambda) = \frac{1}{2} ||x - Dz||_2^2 - \lambda(t - ||z||_1),$$

and the subgradient with respect to z is

$$\partial_z \mathcal{L}(z,\lambda) = -D^T(x - Dz) + \lambda v,$$

where $v_j = 1$ if $z_j > 0$, $v_j = -1$ if $z_j < 0$, and $v_j \in [-1, 1]$ if $z_j = 0$. From the definition of v, it follows that

$$v^T z = \|z\|_1.$$

At an optimal point α , $\partial_z \mathcal{L}(\alpha, \lambda) = 0$, and hence

$$D^{T}(x - D\alpha) = \lambda v$$

$$\updownarrow$$

$$(x - D\alpha)^{T}D = \lambda v^{T}$$

$$\Downarrow$$

$$(x - D\alpha)^{T}D\alpha = \lambda v^{T}\alpha$$

$$\updownarrow$$

$$(x - D\alpha)^{T}D\alpha = \lambda \|\alpha\|_{1},$$

as claimed.

Now, we use the fact that the values of the optimal solutions are close (Lemma 10):

$$\left|v_D^* - v_{\tilde{D}}^*\right| \le \frac{\varepsilon}{\lambda}.$$

But v_D^* is just

$$\begin{split} \frac{1}{2}\langle x - D\alpha, x - D\alpha \rangle + \lambda \|\alpha\|_1 &= \frac{1}{2}\langle x - D\alpha, x - D\alpha \rangle + \langle x - D\alpha, D\alpha \rangle \\ &= \frac{1}{2}\langle x, x - D\alpha \rangle - \frac{1}{2}\langle x - D\alpha, D\alpha \rangle + \langle x - D\alpha, D\alpha \rangle \\ &= \frac{1}{2}\left(\langle x, x - D\alpha \rangle + \langle x - D\alpha, D\alpha \rangle\right) \\ &= \frac{1}{2}\langle x + D\alpha, x - D\alpha \rangle \\ &= \frac{1}{2}\left(\|x\|_2^2 - \|D\alpha\|_2^2\right). \end{split}$$

Consequently,

$$\left| \frac{1}{2} \left(\|x\|_2^2 - \|D\alpha\|_2^2 \right) - \frac{1}{2} \left(\|x\|_2^2 - \|\tilde{D}\tilde{\alpha}\|_2^2 \right) \right| \le \frac{\varepsilon}{\lambda}$$

and hence

$$\left| \|D\alpha\|_2^2 - \|\tilde{D}\tilde{\alpha}\|_2^2 \right| \le \frac{2\varepsilon}{\lambda}.$$

Finally, we prove stability of the optimal reconstructor. Rather than showing that $||D\alpha - \tilde{D}\tilde{\alpha}||_2^2$ is $O(\varepsilon)$, it will be more convenient for later purposes to prove the following roughly equivalent result.

Lemma 12 (Reconstructor Stability) If $||D - \tilde{D}||_2 \le \varepsilon$, then

$$||D\alpha - D\tilde{\alpha}||_2^2 \le \frac{26\varepsilon}{\lambda}.$$

Proof: Let $\alpha' := \frac{1}{2}(\alpha + \tilde{\alpha})$. From the optimality of α , it follows that $v_D(\alpha) \leq v_D(\alpha')$, or more explicitly:

$$\frac{1}{2}\|x - D\alpha\|_2^2 + \lambda \|\alpha\|_1 \le \frac{1}{2}\|x - D\alpha'\|_2^2 + \lambda \|\alpha'\|_1.$$
(18)

First, note that $\left| \|D\tilde{\alpha}\|_2^2 - \|\tilde{D}\tilde{\alpha}\|_2^2 \right| \leq \frac{5\varepsilon}{\lambda}$, because

$$\begin{split} \left| \|D\tilde{\alpha}\|_{2}^{2} - \|\tilde{D}\tilde{\alpha}\|_{2}^{2} \right| &\leq 2 \left| \langle D\tilde{\alpha}, (\tilde{D} - D)\alpha \rangle \right| + \|(\tilde{D} - D)\tilde{\alpha}\|_{2}^{2} \\ &\leq 2 \|D\tilde{\alpha}\|_{2} \|\tilde{D} - D\|_{2} \|\tilde{\alpha}\|_{2} + \left(\|\tilde{D} - D\|_{2} \|\tilde{\alpha}\|_{2} \right)^{2} \\ &\leq 2 \left(1 + \frac{\varepsilon}{\lambda} \right) \frac{\varepsilon}{\lambda} + \left(\frac{\varepsilon}{\lambda} \right)^{2} \\ &\leq \frac{5\varepsilon}{\lambda}, \end{split}$$

assuming $\varepsilon \leq \lambda$. Combining this fact with Lemma 11, $\left| \|D\alpha\|_2^2 - \|\tilde{D}\tilde{\alpha}\|_2^2 \right| \leq \frac{2\varepsilon}{\lambda}$, yields

$$\left| \|D\alpha\|_2^2 - \|D\tilde{\alpha}\|_2^2 \right| \le \frac{7\varepsilon}{\lambda}.$$

By the convexity of the 1-norm, the RHS of (18) obeys:

$$\begin{split} &\frac{1}{2} \left\| x - D\left(\frac{\alpha + \tilde{\alpha}}{2}\right) \right\|_{2}^{2} + \lambda \left\| \frac{\alpha + \tilde{\alpha}}{2} \right\|_{1} \\ &\leq \frac{1}{2} \left\| x - \frac{1}{2} (D\alpha + D\tilde{\alpha}) \right\|_{2}^{2} + \frac{\lambda}{2} \|\alpha\|_{1} + \frac{\lambda}{2} \|\tilde{\alpha}\|_{1} \\ &= \frac{1}{2} \left(\|x\|_{2}^{2} - 2\langle x, \frac{1}{2} (D\alpha + D\tilde{\alpha}) \rangle + \frac{1}{4} \|D\alpha + D\tilde{\alpha}\|_{2}^{2} \right) + \frac{\lambda}{2} \|\alpha\|_{1} + \frac{\lambda}{2} \|\tilde{\alpha}\|_{1} \\ &= \frac{1}{2} \|x\|_{2}^{2} - \frac{1}{2} \langle x, D\alpha \rangle - \frac{1}{2} \langle x, D\tilde{\alpha} \rangle + \frac{1}{8} \left(\|D\alpha\|_{2}^{2} + \|D\tilde{\alpha}\|_{2}^{2} + 2\langle D\alpha, D\tilde{\alpha} \rangle \right) + \frac{\lambda}{2} \|\alpha\|_{1} + \frac{\lambda}{2} \|\tilde{\alpha}\|_{1} \\ &\leq \frac{1}{2} \|x\|_{2}^{2} - \frac{1}{2} \langle x, D\alpha \rangle - \frac{1}{2} \langle x, D\tilde{\alpha} \rangle + \frac{1}{4} \|D\alpha\|_{2}^{2} + \frac{1}{4} \langle D\alpha, D\tilde{\alpha} \rangle + \frac{\lambda}{2} \|\alpha\|_{1} + \frac{\lambda}{2} \|\tilde{\alpha}\|_{1} + \frac{7}{8} \frac{\varepsilon}{\lambda} \\ &\leq \frac{1}{2} \|x\|_{2}^{2} - \frac{1}{2} \langle x, D\alpha \rangle - \frac{1}{2} \langle x, D\tilde{\alpha} \rangle + \frac{1}{4} \|D\alpha\|_{2}^{2} + \frac{1}{4} \langle D\alpha, D\tilde{\alpha} \rangle + \frac{1}{2} \langle x - D\alpha, D\alpha \rangle + \frac{1}{2} \langle x - \tilde{D}\tilde{\alpha}, \tilde{D}\tilde{\alpha} \rangle + \frac{7}{8} \frac{\varepsilon}{\lambda} \\ &\leq \frac{1}{2} \|x\|_{2}^{2} - \frac{1}{2} \langle x, D\alpha \rangle - \frac{1}{2} \langle x, D\tilde{\alpha} \rangle + \frac{1}{4} \|D\alpha\|_{2}^{2} + \frac{1}{4} \langle D\alpha, D\tilde{\alpha} \rangle + \frac{1}{2} \langle x, D\alpha \rangle - \frac{1}{2} \|D\alpha\|_{2}^{2} + \frac{1}{2} \langle x, D\tilde{\alpha} \rangle - \frac{1}{2} \|D\alpha\|_{2}^{2} \\ &+ \left(\frac{7}{8} + \frac{3}{2}\right) \frac{\varepsilon}{\lambda} \end{split}$$

which simplifies to

$$\begin{split} &\frac{1}{2}\|x\|_2^2 - \frac{3}{4}\|D\alpha\|_2^2 - \frac{1}{2}\langle x, D\alpha\rangle - \frac{1}{2}\langle x, D\tilde{\alpha}\rangle + \frac{1}{4}\langle D\alpha, D\tilde{\alpha}\rangle + \frac{1}{2}\langle x, D\alpha\rangle + \frac{1}{2}\langle x, D\tilde{\alpha}\rangle + \frac{19}{8}\frac{\varepsilon}{\lambda} \\ &= \frac{1}{2}\|x\|_2^2 - \frac{3}{4}\|D\alpha\|_2^2 + \frac{1}{4}\langle D\alpha, D\tilde{\alpha}\rangle + \frac{19}{8}\frac{\varepsilon}{\lambda}. \end{split}$$

Now, taking the (expanded) LHS of (18) and the newly derived upper bound of the RHS of (18) yields the inequality:

$$\begin{split} &\frac{1}{2}\|x\|_{2}^{2} - \langle x, D\alpha \rangle + \frac{1}{2}\|D\alpha\|_{2}^{2} + \lambda\|\alpha\|_{1} \\ &\leq \frac{1}{2}\|x\|_{2}^{2} - \frac{3}{4}\|D\alpha\|_{2}^{2} + \frac{1}{4}\langle D\alpha, D\tilde{\alpha} \rangle + \frac{19}{8}\frac{\varepsilon}{\lambda}. \end{split}$$

which implies that

$$-\langle x, D\alpha \rangle + \frac{1}{2} \|D\alpha\|_{2}^{2} + \lambda \|\alpha\|_{1}$$

$$\leq -\frac{3}{4} \|D\alpha\|_{2}^{2} + \frac{1}{4} \langle D\alpha, D\tilde{\alpha} \rangle + \frac{19}{8} \frac{\varepsilon}{\lambda}.$$

Replacing $\lambda \|\alpha\|_1$ with $\langle x - D\alpha, D\alpha \rangle$ yields:

$$-\langle x, D\alpha \rangle + \frac{1}{2} \|D\alpha\|_{2}^{2} + \langle x, D\alpha \rangle - \|D\alpha\|_{2}^{2}$$

$$\leq -\frac{3}{4} \|D\alpha\|_{2}^{2} + \frac{1}{4} \langle D\alpha, D\tilde{\alpha} \rangle + \frac{19}{8} \frac{\varepsilon}{\lambda},$$

implying that

$$\frac{1}{4} \|D\alpha\|_2^2 \le \frac{1}{4} \langle D\alpha, D\tilde{\alpha} \rangle + \frac{19}{8} \frac{\varepsilon}{\lambda}.$$

Hence,

$$||D\alpha||_2^2 \le \langle D\alpha, D\tilde{\alpha} \rangle + \frac{19}{2} \frac{\varepsilon}{\lambda}.$$

Now, note that

$$\begin{split} \|D\alpha - D\tilde{\alpha}\|_2^2 &= \|D\alpha\|_2^2 + \|D\tilde{\alpha}\|_2^2 - 2\langle D\alpha, D\tilde{\alpha}\rangle \\ &\leq \|D\alpha\|_2^2 + \|D\tilde{\alpha}\|_2^2 - 2\|D\alpha\|_2^2 + 19\frac{\varepsilon}{\lambda} \\ &\leq \|D\alpha\|_2^2 + \|D\alpha\|_2^2 - 2\|D\alpha\|_2^2 + 26\frac{\varepsilon}{\lambda} \\ &= 26\frac{\varepsilon}{\lambda}. \end{split}$$

A.3 The sparsity lemma

We now prove that the solution to the perturbed problem is sparse for sufficiently small ε .

Lemma 13 (Preservation of Sparsity) Under Assumptions (A1)-(A3), if

$$\tau \ge \varepsilon \left(1 + \frac{1}{\lambda} \right) + \sqrt{\frac{26\varepsilon}{\lambda}},$$

then $\tilde{\alpha}_i = 0$ for all $i \in \mathcal{I}$.

Proof: Let $i \in \mathcal{I}$ be arbitrary. To prove that $\tilde{\alpha}_i = 0$, it is sufficient to show that

$$\left| \langle \tilde{D}_i, x - \tilde{D}\tilde{\alpha} \rangle \right| < \lambda,$$

since $\tilde{\alpha}_i$ is hence zero.

First, note that

$$\begin{aligned} \left| \langle \tilde{D}_{i}, x - \tilde{D}\tilde{\alpha} \rangle \right| &= \left| \langle D_{i} + \tilde{D}_{i} - D_{i}, x - \tilde{D}\tilde{\alpha} \rangle \right| \\ &\leq \left| \langle D_{i}, x - \tilde{D}\tilde{\alpha} \rangle \right| + \|\tilde{D}_{i} - D_{i}\|_{2} \|x - \tilde{D}\tilde{\alpha}\|_{2} \\ &\leq \left| \langle D_{i}, x - \tilde{D}\tilde{\alpha} \rangle \right| + \varepsilon \qquad \text{(since } \|x\|_{2} \leq 1) \end{aligned}$$

and

$$\begin{split} \left| \langle D_i, x - \tilde{D}\tilde{\alpha} \rangle \right| &= \left| \langle D_i, x - (D + \tilde{D} - D)\tilde{\alpha} \rangle \right| \\ &\leq \left| \langle D_i, x - D\tilde{\alpha} \rangle \right| + \left| \langle D_i, (\tilde{D} - D)\tilde{\alpha} \rangle \right| \\ &\leq \left| \langle D_i, x - D\tilde{\alpha} \rangle \right| + \|D_i\|_2 \|\tilde{D} - D\|_2 \|\tilde{\alpha}\|_2 \\ &\leq \left| \langle D_i, x - D\tilde{\alpha} \rangle \right| + \frac{\varepsilon}{\lambda}. \end{split}$$

Hence,

$$\left| \langle \tilde{D}_i, x - \tilde{D}\alpha \rangle \right| \le \left| \langle D_i, x - D\tilde{\alpha} \rangle \right| + \varepsilon \left(1 + \frac{1}{\lambda} \right),$$

and so it is sufficient to show that

$$|\langle D_i, x - D\tilde{\alpha} \rangle| < \lambda - \varepsilon \left(1 + \frac{1}{\lambda}\right).$$

Now,

$$|\langle D_{i}, x - D\tilde{\alpha} \rangle| = |\langle D_{i}, x - D\tilde{\alpha} + D\alpha - D\alpha \rangle|$$

$$\leq |\langle D_{i}, x - D\alpha \rangle| + |\langle D_{i}, D\alpha - D\tilde{\alpha} \rangle|$$

$$< \lambda - \tau + ||D_{i}||_{2} ||D\alpha - D\tilde{\alpha}||_{2}$$

$$< \lambda - \tau + \sqrt{\frac{26\varepsilon}{\lambda}}, \tag{19}$$

where (19) is due to Lemma 12. Consequently, it is sufficient if τ is chosen to satisfy

$$\lambda - \tau + \sqrt{\frac{26\varepsilon}{\lambda}} \le \lambda - \varepsilon \left(1 + \frac{1}{\lambda}\right),$$

yielding:

$$\tau \ge \varepsilon \left(1 + \frac{1}{\lambda} \right) + \sqrt{\frac{26\varepsilon}{\lambda}}.$$

A.4 Proof of Sparse Coding Stability Theorem

Proof (of Theorem 1): Recall that $\varphi_D(x)$ is the unique optimal solution to the problem

$$\min_{z \in \mathbb{R}^k} \frac{1}{2} ||x - Dz||_2^2 + \lambda_1 ||z||_1.$$

If not for ℓ_1 penalty, in standard form, the quadratic program is

$$\min_{z \in \mathbb{R}^k} z^T D^T D z - z^T (2Dx) + \lambda_1 ||x||_1$$

Similarly, let $\tilde{Q}(\cdot)$ be the objective using \tilde{D} instead of D. Denoting $\bar{z}:=\begin{pmatrix}z\\z^+\\z^-\end{pmatrix}$, an equivalent formulation is

For optimal solutions $\bar{z}_* := \begin{pmatrix} z_* \\ z_*^+ \\ z_*^- \end{pmatrix}$ and $\bar{t}_* := \begin{pmatrix} t_* \\ t_*^+ \\ t_*^- \end{pmatrix}$ of Q and \tilde{Q} respectively, from Daniel (1973), we have

$$(\bar{u} - \bar{z}_*)^T \nabla Q(\bar{z}_*) \ge 0 \tag{20}$$

$$(\bar{u} - \bar{t}_*)^T \nabla \tilde{Q}(\bar{t}_*) \ge 0 \tag{21}$$

for all $\bar{u} \in \mathbb{R}^{3k}$. Setting \bar{u} to \bar{t}_* in (20) and \bar{u} to \bar{z}_* in (21) and adding (21) and (20) yields

$$(\bar{t}_* - \bar{z}_*)^T (\nabla Q(\bar{z}_*) - \nabla \tilde{Q}(\bar{t}_*)) \ge 0,$$

which is equivalent to

$$(\bar{t}_* - \bar{z}_*)^T (\nabla \tilde{Q}(\bar{t}_*) - \nabla \tilde{Q}(\bar{z}_*)) \le (\bar{t}_* - \bar{z}_*)^T (\nabla Q(\bar{z}_*) - \nabla \tilde{Q}(\bar{z}_*))$$
(22)

Here,

$$\nabla Q(z) = \frac{1}{2} \begin{pmatrix} D^T D & \mathbf{0}_{k \times 2k} \\ \mathbf{0}_{2k \times k} & \mathbf{0}_{2k \times 2k} \end{pmatrix} z - \frac{1}{2} \begin{pmatrix} 2D^T \\ \mathbf{0}_{2k \times d} \end{pmatrix} x + \lambda_1 \begin{pmatrix} \mathbf{0}_k \\ \mathbf{1}_{2k} \end{pmatrix}.$$

After plugging in the expansions of ∇Q and $\nabla \tilde{Q}$ and incurring cancellations from the zeros, (22) becomes

$$(t_{*}-z_{*})^{T}\tilde{D}^{T}\tilde{D}(t_{*}-z_{*}) \leq (t_{*}-z_{*})^{T}\left((D^{T}D-\tilde{D}^{T}\tilde{D})z_{*}+2(\tilde{D}-D)^{T}x\right)$$

$$\leq (t_{*}-z_{*})^{T}(D^{T}D-\tilde{D}^{T}\tilde{D})z_{*}+2\|t_{*}-z_{*}\|_{2}\|(\tilde{D}-D)^{T}x\|_{2}$$

$$\leq (t_{*}-z_{*})^{T}(D^{T}D-\tilde{D}^{T}\tilde{D})z_{*}+\|t_{*}-z_{*}\|_{2}(2\varepsilon)$$

$$(24)$$

Let us gain a handle on the first term. Note that $\tilde{D} = D + E$ for some E satisfying $||E||_2 \le \varepsilon$. Hence,

$$\begin{split} &(t_{*}-z_{*})^{T}(D^{T}D-\tilde{D}^{T}\tilde{D})z_{*}\\ &=\left|(t_{*}-z_{*})^{T}(E^{T}D+D^{T}E+E^{T}E)z_{*}\right|\\ &\leq\left|(t_{*}-z_{*})^{T}E^{T}Dz_{*}\right|+\left|(t_{*}-z_{*})^{T}D^{T}Ez_{*}\right|+\left|(t_{*}-z_{*})^{T}E^{T}Ez_{*}\right|\\ &\leq\left\|t_{*}-z_{*}\right\|_{2}\left(\|E\|_{2}\|Dz_{*}\|_{2}+\|D(t_{*}-z_{*})\|_{2}\|Ez_{*}\|_{2}+\|t_{*}-z_{*}\|_{2}\|E\|_{2}^{2}\|z_{*}\|_{2}\right)\\ &\leq\left\|t_{*}-z_{*}\right\|_{2}\left(\frac{\varepsilon\sqrt{s}}{\lambda}+\frac{\varepsilon\sqrt{s}}{\lambda}+\frac{\varepsilon^{2}}{\lambda}\right)\\ &\leq\left\|t_{*}-z_{*}\right\|_{2}\frac{3\varepsilon\sqrt{s}}{\lambda}, \end{split}$$

where the last step follows because if $||z_*||_0 \le s$, then Lemma 14 in Appendix E implies that $||Dz_*||_2 \le \sqrt{s}||z_*||_2$ (and $||z_*||_2 \le ||z_*||_1 \le \frac{1}{\lambda}$).

Now, observe from Lemma 13 that $||t_* - z_*||_0 \le s$. Combining this result with the fact that \tilde{D} has s-incoherence lower bounded by μ implies the desired result:

$$||t_* - z_*||_2 \le \frac{3\varepsilon\sqrt{s}}{\lambda\mu}.$$

B Proof of Restricted Stability Theorem

Proof (of Theorem 2): For convenience, define $\mathcal{A} := \operatorname{supp} \varphi_D(x)$, let α be equal to $(\varphi_D(x))_{\mathcal{A}}$, and define the scaled sign vector $\zeta := \lambda \operatorname{sign}(\alpha)$. Our strategy will be to show that, for some $\Delta \in \mathbb{R}^s$, the optimal perturbed solution $\varphi_{\tilde{D}}(x)$ satisfies $(\varphi_{\tilde{D}}(x))_{\mathcal{A}} = \alpha + \Delta$ and $(\varphi_{\tilde{D}}(x))_{\mathcal{A}^c} = 0$, where $\mathcal{A}^c := [k] \setminus \mathcal{A}$.

From the optimality conditions for the LASSO (e.g. see optimality conditions L1 and L2 of Asif and Romberg (2010)), it is sufficient to find Δ such that

$$\langle \tilde{D}_j, x - \tilde{D}_{\mathcal{A}}(\alpha + \Delta) \rangle = \zeta_j \quad \text{if } j \in \mathcal{A},$$

 $\left| \langle \tilde{D}_j, x - \tilde{D}_{\mathcal{A}}(\alpha + \Delta) \rangle \right| < \lambda \quad \text{otherwise.}$

We proceed by setting up the linear system and characterizing the solution vector Δ :

$$\tilde{D}_{\mathcal{A}}^{T}(x - \tilde{D}_{\mathcal{A}}(\alpha + \Delta)) = \zeta \xrightarrow{\text{Solve for } \Delta} \Delta = (\tilde{D}_{\mathcal{A}}^{T}\tilde{D}_{\mathcal{A}})^{-1}(\tilde{D}_{\mathcal{A}}^{T}(x - \tilde{D}_{\mathcal{A}}\alpha) - \zeta).$$

Since $\tilde{D} = D + E$ for $||E||_2 \le \varepsilon$,

$$\begin{split} \tilde{D}_{\mathcal{A}}^T(x - \tilde{D}_{\mathcal{A}}\alpha) &= (D_{\mathcal{A}} + E_{\mathcal{A}})^T(x - (D_{\mathcal{A}} + E_{\mathcal{A}})\alpha) \\ &= D_{\mathcal{A}}^T(x - D_{\mathcal{A}}\alpha) - D_{\mathcal{A}}^T E_{\mathcal{A}}\alpha + E_{\mathcal{A}}^T(x - (D_{\mathcal{A}} + E_{\mathcal{A}})\alpha) \\ &= \zeta - D_{\mathcal{A}}^T E_{\mathcal{A}}\alpha + E_{\mathcal{A}}^T(x - (D_{\mathcal{A}} + E_{\mathcal{A}})\alpha), \end{split}$$

and so the solution for Δ can be reformulated as

$$\Delta = (\tilde{D}_{\mathcal{A}}^T \tilde{D}_{\mathcal{A}})^{-1} (-D_{\mathcal{A}}^T E_{\mathcal{A}} \alpha + E_{\mathcal{A}}^T (x - (D_{\mathcal{A}} + E_{\mathcal{A}}) \alpha)).$$

Now,

$$\begin{split} \|\Delta\|_{2} &\leq \|\tilde{D}_{\mathcal{A}}^{T} \tilde{D}_{\mathcal{A}})^{-1} \|_{2} (\|D_{\mathcal{A}}^{T} E_{\mathcal{A}} \alpha\|_{2} + \|E_{\mathcal{A}}^{T} (x - (D_{\mathcal{A}} + E_{\mathcal{A}}) \alpha)\|_{2}) \\ &\leq \frac{1}{\mu} (\frac{\varepsilon \sqrt{s}}{\lambda} + \varepsilon) \\ &= \frac{\varepsilon}{\mu} \left(\frac{\sqrt{s}}{\lambda} + 1\right). \end{split}$$

For $y \in \mathbb{R}^s$, let $y_{k \times 1}$ be the extension to \mathbb{R}^k satisfying $(y_{k \times 1})_{\mathcal{A}} = y$ and $(y_{k \times 1})_{\mathcal{A}^c} = 0$. For $(\alpha + \Delta)_{k \times 1}$ to be optimal for LASSO (λ, \tilde{D}, x) , $(\alpha + \Delta)_{k \times 1}$ must satisfy the two optimality conditions and Δ must be small enough such that sign consistency holds between α and $(\alpha + \Delta)$ (i.e. $\operatorname{sign}(\alpha_j) = \operatorname{sign}(\alpha_j + \Delta_j)$ for all $j \in [s]$).

We first check the optimality conditions. The first optimality condition is equivalent to

$$\langle \tilde{D}_j, x - \tilde{D}_{\mathcal{A}}(\alpha + \Delta) \rangle = \lambda \quad \text{for } j \in \mathcal{A};$$

this condition is satisfied by construction. The second optimality condition is equivalent to

$$\left| \langle \tilde{D}_j, x - \tilde{D}_{\mathcal{A}}(\alpha + \Delta) \rangle \right| < \lambda \quad \text{ for } j \notin \mathcal{A}.$$

But for $j \notin \mathcal{A}$,

$$\begin{split} \left| \langle \tilde{D}_{j}, x - \tilde{D}_{\mathcal{A}}(\alpha + \Delta) \rangle \right| &= \left| \langle D_{j} + E_{j}, x - (D_{\mathcal{A}} + E_{\mathcal{A}})(\alpha + \Delta) \rangle \right| \\ &= \left| \langle D_{j}, x - D_{\mathcal{A}} \alpha \rangle - \langle D_{j}, D_{\mathcal{A}} \Delta \rangle + \langle E_{j}, x - (D_{\mathcal{A}} + E_{\mathcal{A}})(\alpha + \Delta) \rangle - \langle D_{j}, E_{\mathcal{A}}(\alpha + \Delta) \rangle \right| \\ &< \lambda - \tau + \frac{\varepsilon \sqrt{s}}{\mu} \left(\frac{\sqrt{s}}{\lambda} + 1 \right) + \varepsilon + \frac{\varepsilon}{\lambda} \\ &= \lambda - \tau + \varepsilon \left(\frac{\frac{s}{\lambda} + \sqrt{s}}{\mu} + \frac{1}{\lambda} + 1 \right), \end{split}$$

and so this condition is satisfied provided that

$$\varepsilon \left(\frac{\frac{s}{\lambda} + \sqrt{s}}{\mu} + \frac{1}{\lambda} + 1 \right) \le \tau.$$

Now, we check sign consistency. Clearly sign consistency holds over \mathcal{A}^c . It remains to check that it holds over \mathcal{A} . Observe that

$$\|\Delta\|_{\infty} \le \|\Delta\|_2 \le \frac{\varepsilon}{\mu} \left(\frac{\sqrt{s}}{\lambda} + 1\right).$$

Hence, sign consistency holds provided that

$$|\alpha_i| > \varepsilon \left(\frac{1}{\mu} \left(\frac{\sqrt{s}}{\lambda} + 1\right)\right).$$

All the above constraints are satisfied if τ satisfies

$$\varepsilon \left(\frac{\frac{s}{\lambda} + \sqrt{s}}{\mu} + \frac{1}{\lambda} + 1 \right) \le \tau.$$

C Proof of Symmetrization by Ghost Sample Lemma

Proof (of Lemma 1): Replace $\mathcal{F}(\sigma_n)$ from the notation of Mendelson and Philips (2004) with $\mathcal{F}(\mathbf{z}, \mathbf{x}'')$. A modified one-sided version of (Mendelson and Philips, 2004, Lemma 2.2) that uses the more favorable Chebyshev-Cantelli inequality implies that, for every t > 0:

$$\left(1 - \frac{4 \sup_{f \in \mathcal{F}} \operatorname{Var}(l(\cdot, f))}{4 \sup_{f \in \mathcal{F}} \operatorname{Var}(l(\cdot, f)) + mt^{2}}\right) \operatorname{Pr}_{\mathbf{z}, \mathbf{x}''} \left\{ \exists f \in \mathcal{F}(\mathbf{z}, \mathbf{x}''), \ (P - P_{\mathbf{z}}) l(\cdot, f) \ge t \right\}
\leq \operatorname{Pr}_{\mathbf{z}, \mathbf{z}', \mathbf{x}''} \left\{ \exists f \in \mathcal{F}(\mathbf{z}, \mathbf{x}''), \ (P_{\mathbf{z}'} - P_{\mathbf{z}}) l(\cdot, f) \ge \frac{t}{2} \right\}.$$

As the losses lie in [0, b] by assumption, it follows that $\sup_{f \in \mathcal{F}} \operatorname{Var}(l(\cdot, f)) \leq \frac{b^2}{4}$. The lemma follows since the left hand factor of the LHS of the above inequality is at least $\frac{1}{2}$ whenever $m \geq \left(\frac{b}{t}\right)^2$.

D Proofs for overcomplete setting

Proof (of Theorem 3): Proposition 1 and Lemmas 5 and 6 imply that

$$\begin{split} & \Pr_{\mathbf{z}} \left\{ \begin{array}{l} \exists f \in \mathcal{F}_{\mu}, \ \left[\operatorname{margin}_{s}(D, \mathbf{x}) > \iota \right] \\ & \quad \text{and} \ \left((P - P_{\mathbf{z}})l(\cdot, f) > t \right) \end{array} \right\} \\ & \leq 2 \left(\left(\frac{8(r/2)^{1/(d+1)}}{\varepsilon} \right)^{(d+1)k} \exp(-m\varpi^{2}/(2b^{2})) + \delta \right). \end{split}$$

Equivalently,

$$\begin{split} \Pr_{\mathbf{z}} \left\{ \begin{array}{l} \exists f \in \mathcal{F}_{\mu}, \ \left[\operatorname{margin}_{s}(D, \mathbf{x}) > \iota \right] \\ \text{and} \ \left((P - P_{\mathbf{z}}) l(\cdot, f) > 2 \left(\varpi + 2L\beta + \frac{b\eta(m, d, k, \varepsilon, \delta)}{m} \right) \right) \end{array} \right\} \\ \leq 2 \left(\left(\frac{8(r/2)^{1/(d+1)}}{\varepsilon} \right)^{(d+1)k} \exp(-m\varpi^{2}/(2b^{2})) + \delta \right). \end{split}$$

Now, expand β and η and replace δ with $\delta/4$:

$$\begin{split} \Pr_{\mathbf{z}} \left\{ \begin{array}{l} \exists f \in \mathcal{F}_{\mu}, \ \left[\operatorname{margin}_{s}(D, \mathbf{x}) > \iota \right] \ \mathbf{and} \\ (P - P_{\mathbf{z}}) l(\cdot, f) > 2 \left(\varpi + 2L\varepsilon \frac{1}{\lambda} \left(1 + \frac{3r\sqrt{s}}{\mu} \right) + \frac{b(dk \log \frac{3096}{\operatorname{margin}_{s}^{2}(D, \mathbf{x}) \cdot \lambda} + \log(2m+1) + \log \frac{4}{\delta})}{m} \right) \end{array} \right\} \\ \leq 2 \left(\frac{8(r/2)^{1/(d+1)}}{\varepsilon} \right)^{(d+1)k} \exp(-m\varpi^{2}/(2b^{2})) + \frac{\delta}{2}. \end{split}$$

Choosing
$$\frac{\delta}{4} = \left(\frac{8(r/2)^{1/(d+1)}}{\varepsilon}\right)^{(d+1)k} \exp(-m\varpi^2/(2b^2))$$
 yields

$$\Pr_{\mathbf{z}} \left\{ \begin{array}{l} \exists f \in \mathcal{F}_{\mu}, \ \left[\operatorname{margin}_{s}(D, \mathbf{x}) > \iota \right] \ \mathbf{and} \\ (P - P_{\mathbf{z}}) l(\cdot, f) > 2 \left(\begin{array}{l} \varpi + 2L\varepsilon \frac{1}{\lambda} \left(1 + \frac{3r\sqrt{s}}{\mu} \right) + \\ \frac{b(dk \log \frac{3096}{\operatorname{margin}_{s}^{2}(D, \mathbf{x}) \cdot \lambda} + \log(2m+1) + (d+1)k \log \frac{\varepsilon}{8(r/2)^{1/(d+1)}} + \frac{m\varpi^{2}}{b^{2}})}{m} \end{array} \right) \right\} \\ \leq 4 \cdot \left(\frac{8(r/2)^{1/(d+1)}}{\varepsilon} \right)^{(d+1)k} \exp(-m\varpi^{2}/(2b^{2})), \end{array}$$

which is equivalent to

$$\Pr_{\mathbf{z}} \left\{ \begin{array}{l} \exists f \in \mathcal{F}_{\mu}, \ \left[\operatorname{margin}_{s}(D, \mathbf{x}) > \iota \right] \ \mathbf{and} \\ (P - P_{\mathbf{z}})l(\cdot, f) > 2 \left(\begin{array}{l} \varpi + 2L\varepsilon \frac{1}{\lambda} \left(1 + \frac{3r\sqrt{s}}{\mu} \right) + \\ \frac{b(dk \log \frac{3096}{\operatorname{margin}_{s}^{2}(D, \mathbf{x}) \cdot \lambda} - (d+1)k \log \frac{8}{\varepsilon} + k \log \frac{2}{r} + \log(2m+1) + \frac{m\varpi^{2}}{b^{2}})}{m} \end{array} \right) \right. \right\} \\ \leq 4 \cdot \left(\frac{8(r/2)^{1/(d+1)}}{\varepsilon} \right)^{(d+1)k} \exp(-m\varpi^{2}/(2b^{2})), \end{array}$$

Let δ (a new variable, not related to the previous incarnation of δ) be equal to the upper bound, and solve for ϖ , yielding:

$$\varpi = b\sqrt{\frac{2((d+1)k\log\frac{8}{\varepsilon} + k\log\frac{r}{2} + \log\frac{4}{\delta})}{m}}$$

and hence

$$\Pr_{\mathbf{z}} \left\{ \begin{array}{l} \exists f = (D, w) \in \mathcal{F}_{\mu}, \; \left[\operatorname{margin}_{s}(D, \mathbf{x}) > \iota \right] \; \mathbf{and} \\ (P - P_{\mathbf{z}})l(\cdot, f) > 2 \left(\begin{array}{l} b\sqrt{\frac{2((d+1)k\log\frac{8}{\delta} + k\log\frac{r}{2} + \log\frac{4}{\delta})}{m}} + 2L\varepsilon\frac{1}{\lambda}\left(1 + \frac{3r\sqrt{s}}{\mu}\right) + \\ \frac{b(dk\log\frac{3096}{\operatorname{margin}_{s}^{2}(D, \mathbf{x}) \cdot \lambda} + \log(2m+1) + \log\frac{4}{\delta})}{m} \end{array} \right) \right\} \\ < \delta. \end{array} \right\}$$

If we set $\varepsilon = \frac{1}{m}$, then provided that $m > \frac{387}{\text{margin}_s^2(D, \mathbf{x}) \cdot \lambda}$:

$$\Pr_{\mathbf{z}} \left\{ \begin{array}{l} \exists f \in \mathcal{F}_{\mu}, \ \left[\operatorname{margin}_{s}(D, \mathbf{x}) > \iota \right] \ \mathbf{and} \\ (P - P_{\mathbf{z}})l(\cdot, f) > 2 \left(\begin{array}{l} b\sqrt{\frac{2((d+1)k\log(8m) + k\log\frac{r}{2} + \log\frac{4}{\delta})}{m}} + \frac{2L}{m} \left(\frac{1}{\lambda}(1 + \frac{3r\sqrt{s}}{\mu}) \right) + \\ \frac{b}{m} \left(dk\log\frac{3096}{\operatorname{margin}_{s}^{2}(D, \mathbf{x}) \cdot \lambda} + \log(2m + 1) + \log\frac{4}{\delta} \right) \end{array} \right) \right. \right\} \\ \leq \delta. \end{array}$$

It remains to distribute a prior across the bounds for each choice of s and μ . To each choice of $s \in [k]$ assign prior probability $\frac{1}{k}$. To each choice of $i \in \mathbb{N} \cup \{0\}$ for $2^{-i} \leq \mu$ assign prior probability $(i+1)^{-2}$. For a given choice of $s \in [k]$ and $2^{-i} \leq \mu$ we use $\delta(s,i) := \frac{6}{\pi^2} \frac{1}{(i+1)^2} \frac{1}{k} \delta$ (since $\sum_{i=1}^{\infty} \frac{1}{i^2} = \frac{\pi^2}{6}$). Then, provided that

$$m > \frac{387}{\text{margin}_s(D, \mathbf{x})^2 \lambda},$$

the generalization error $(P - P_{\mathbf{z}})l(\cdot, f)$ is bounded by:

$$2b\sqrt{\frac{2\left((d+1)k\log(8m) + k\log\frac{r}{2} + \log\frac{2\pi^{2}\left(\log_{2}\frac{4}{\mu_{s}(D)}\right)^{2}k\right)}{m}}{m}} + \frac{2b}{m}\left(dk\log\frac{3096}{\mathrm{margin}_{s}^{2}(D,\mathbf{x})\cdot\lambda} + \log(2m+1) + \log\frac{2\pi^{2}\left(\log_{2}\frac{4}{\mu_{s}(D)}\right)^{2}k}{3\delta}\right) + \frac{4L}{m}\left(\frac{1}{\lambda}(1 + \frac{6r\sqrt{s}}{\mu_{s}(D)})\right).$$

E Infinite-dimensional setting

Proof (of Lemma 7): Recall that $\eta = \log \frac{1}{\delta}$. Suppose, as in the event being measured, that there is no subset of the ghost sample \mathbf{x}' of size at least η such that the τ -level s-margin condition holds for the entire subset. Equivalently, there is a subset of at least η points in the ghost sample \mathbf{x}' that violate the τ -level s-coding margin condition. From the permutation argument, if no point of \mathbf{x}'' violates $\left[\text{margin}_s(D,\cdot) > \tau\right]$, then the probability that over $\eta = \log \frac{1}{\delta}$ points of \mathbf{x}' will violate $\left[\text{margin}_s(D,\cdot) > \tau\right]$ is at most δ .

Proof (of Lemma 9): By definition, $\varphi_{US}(x) = \arg\min_{z \in \mathbb{R}^k} \|x - USv\|_2 + \lambda \|z\|_1$. Note that $\arg\min_{z \in \mathbb{R}^k} \|x - USz\|_2 = \arg\min_{z \in \mathbb{R}^k} \|U^Tx - U^TUSz\|_2 = \arg\min_{z \in \mathbb{R}^k} \|U^Tx - Sz\|_2$, where the first equality follows because any x in the complement of the image of U will be orthogonal to USz, for any z; hence, it is sufficient to approximate the projection of x onto the range of U. Thus, $\varphi_{US}(x) = \arg\min_{z \in \mathbb{R}^k} \|U^Tx - Sz\|_2^2 + \lambda \|z\|_1$. It will be useful to apply a well-known reformulation of this minimization problem as a quadratic program with linear constraints. Denoting $\bar{z} := \bar{z} := (z^T z^{+T} z^{-T})^T$, an equivalent formulation is

For optimal solutions $\bar{z}_* := \begin{pmatrix} z_* \\ z_*^+ \\ z_-^- \end{pmatrix}$ and $\bar{t}_* := \begin{pmatrix} t_* \\ t_*^+ \\ t_-^- \end{pmatrix}$ of Q_U and $Q_{U'}$ respectively, from

Daniel (1973), we have

$$(\bar{u} - \bar{z}_*)^T \nabla Q_U(\bar{z}_*) \ge 0 \tag{25}$$

$$(\bar{u} - \bar{t}_*)^T \nabla Q_{U'}(\bar{t}_*) \ge 0 \tag{26}$$

for all $\bar{u} \in \mathbb{R}^{3k}$. Setting \bar{u} to \bar{t}_* in (25) and \bar{u} to \bar{z}_* in (26) and adding (25) and (26) yields

$$(\bar{t}_* - \bar{z}_*)^T (\nabla Q_U(\bar{z}_*) - \nabla Q_{U'}(\bar{t}_*)) \ge 0,$$

which is equivalent to

$$(\bar{t}_* - \bar{z}_*)^T (\nabla Q_{U'}(\bar{t}_*) - \nabla Q_{U'}(\bar{z}_*)) \le (\bar{t}_* - \bar{z}_*)^T (\nabla Q_U(\bar{z}_*) - \nabla Q_{U'}(\bar{z}_*)). \tag{27}$$

Here, $\nabla Q_U(z) = \begin{pmatrix} S^T S & \mathbf{0}_{k \times 2k} \\ \mathbf{0}_{2k \times k} & \mathbf{0}_{2k \times 2k} \end{pmatrix} z - \begin{pmatrix} 2S^T U^T \\ \mathbf{0}_{2k \times d} \end{pmatrix} x + \lambda \begin{pmatrix} \mathbf{0}_k \\ \mathbf{1}_{2k} \end{pmatrix}$. After plugging in the expansions of ∇Q_U and $\nabla Q_{U'}$ and incurring cancellations from the zeros, (27) becomes

$$(t_* - z_*)^T (S^T S t_* - 2S^T U'^T x - S^T S z_* + 2S^T U'^T x)$$

$$\leq (t_* - z_*)^T (S^T S z_* - 2S^T U^T x - S^T S z_* + 2S^T U'^T x),$$

which reduces to

$$(t_* - z_*)^T S^T S(t_* - z_*) \le 2(t_* - z_*)^T S^T (U'^T - U^T) x.$$

Since both t_* and z_* are s-sparse, wherever we typically would consider the operator norm $||S||_2 := \sup_{|t|=1} ||St||_2$, we instead need only consider the 2s-restricted operator norm $||S||_{2,2s}$.

Note that $(t_* - z_*)^T S^T S(t_* - z_*) \ge \mu_{2s}(S) ||t_* - z_*||_2^2$, which implies that

$$||t_* - z_*||_2^2 \le \frac{2}{\mu_{2s}(S)} ||t_* - z_*|| ||S||_{2,2s} ||(U'^T - U^T)x||$$

and hence

$$||t_* - z_*||_2 \le \frac{2||S||_{2,2s}}{\mu_{2s}(S)} ||(U'^T - U^T)x||_2.$$

Lemma 14 If $S \in (B_{\mathbb{R}^k})^k$, then $||S||_{2,s} \leq \sqrt{s}$.

Proof: Define S_{Λ} as the submatrix of S that selects the columns indexed by Λ . Similarly, for $t \in \mathbb{R}^k$ define the coordinate projection t_{Λ} of t.

$$\begin{split} \sup_{\{t:\|t\|=1,|\operatorname{supp}(t)|\leq s\}} &\|St\|_2 \\ &= \min_{\{\Lambda\subseteq[k]:|\Lambda|\leq s\}} \sup_{\{t:\|t\|=1,\operatorname{supp}(t)\subseteq\Lambda\}} \|S_\Lambda t_\Lambda\|_2 \\ &= \min_{\{\Lambda\subseteq[k]:|\Lambda|\leq s\}} \sup_{\{t:\|t\|=1,\operatorname{supp}(t)\subseteq\Lambda\}} \left\|\sum_{\omega\in\Lambda} t_\omega S_\omega\right\|_2 \\ &\leq \min_{\{\Lambda\subseteq[k]:|\Lambda|\leq s\}} \sup_{\{t:\|t\|=1,\operatorname{supp}(t)\subseteq\Lambda\}} \sum_{\omega\in\Lambda} |t_\omega| \|S_\omega\|_2 \\ &\leq \min_{\{\Lambda\subseteq[k]:|\Lambda|\leq s\}} \sup_{\{t:\|t\|=1,\operatorname{supp}(t)\subseteq\Lambda\}} \sum_{\omega\in\Lambda} |t_\omega| \\ &\leq \min_{\{\Lambda\subseteq[k]:|\Lambda|\leq s\}} \sup_{\{t:\|t\|=1,\operatorname{supp}(t)\subseteq\Lambda\}} |t_\Lambda\|_1 \\ &\leq \min_{\{\Lambda\subseteq[k]:|\Lambda|\leq s\}} \sup_{\{t:\|t\|=1,\operatorname{supp}(t)\subseteq\Lambda\}} \sqrt{s} \|t_\Lambda\|_2 \\ &= \sqrt{s}. \end{split}$$

F Covering numbers

For a Banach space E of dimension d, the ε -covering numbers of the radius r ball of E are bounded as $\mathcal{N}(rB_E,\varepsilon) \leq (4r/\varepsilon)^d$ (Cucker and Smale, 2002, Chapter I, Proposition 5).

For spaces of dictionaries obeying some deterministic property, such as

$$\mathcal{D}_{\mu} = \{ D \in \mathcal{D} : \mu_s(D) \ge \mu \},$$

one must be careful to use a proper ε -cover so that the representative elements of the cover also obey the desired property; a proper cover is more restricted than a cover in that a proper cover must be a subset of the set being covered, rather than simply being a subset of the ambient Banach space. That is, if A is a proper cover of a subset T of a Banach space E, then $A \subseteq T$. For a cover, we need only $A \subseteq E$. The following bound relates proper covering numbers to covering numbers (a simple proof can be found in Vidyasagar 2002, Lemma 2.1): If E is a Banach space and $T \subseteq E$ is a bounded subset, then

$$\mathcal{N}(E, \varepsilon, T) \leq \mathcal{N}_{proper}(E, \varepsilon/2, T).$$

Let $d, k \in \mathbb{N}$. Define $E_{\mu} := \{E \in (B_{\mathbb{R}^d})^k : \mu_s(D) \geq \mu\}$ and $\mathcal{W} := rB_{\mathbb{R}^d}$. The following bounds derive directly from the above.

Proposition 3 The proper ε -covering number of E_{μ} is bounded by $(8/\varepsilon)^{dk}$.

Proposition 4 The product of the proper ε -covering number of E_{μ} and the ε -covering number of W is bounded by

$$\left(\frac{8(r/2)^{1/(d+1)}}{\varepsilon}\right)^{(d+1)k}.$$

References

M. Salman Asif and Justin Romberg. On the LASSO and Dantzig selector equivalence. In *Annual Conference on Information Sciences and Systems (CISS)*, pages 1–6. IEEE, 2010.

David M. Bradley and J. Andrew Bagnell. Differentiable sparse coding. In *Advances in Neural Information Processing Systems 21*, pages 113–120. MIT Press, 2009.

Felipe Cucker and Steve Smale. On the mathematical foundations of learning. Bulletin of the American Mathematical Society, 39:1–49, 2002.

James W. Daniel. Stability of the solution of definite quadratic programs. *Mathematical Program-ming*, 5(1):41–53, 1973.

Maryam Fazel. Matrix rank minimization with applications. *Elec Eng Dept Stanford University*, 54:1–130, 2002.

Sham M. Kakade, Karthik Sridharan, and Ambuj Tewari. On the complexity of linear prediction: Risk bounds, margin bounds, and regularization. In Daphne Koller, Dale Schuurmans, Yoshua Bengio, and Léon Bottou, editors, *Advances in Neural Information Processing Systems 21*, pages 793–800. MIT Press, 2009.

- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Michel Ledoux and Michel Talagrand. Probability in Banach Spaces: isoperimetry and processes. Springer, 1991. ISBN 3540520139.
- Julien Mairal, Francis Bach, Jean Ponce, Guillermo Sapiro, and Andrew Zisserman. Supervised dictionary learning. In Daphne Koller, Dale Schuurmans, Yoshua Bengio, and Léon Bottou, editors, Advances in Neural Information Processing Systems 21, pages 1033–1040. MIT Press, 2009.
- Julien Mairal, Francis Bach, and Jean Ponce. Task-driven dictionary learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(4):791–804, 2012.
- Andreas Maurer and Massimiliano Pontil. K-dimensional coding schemes in hilbert spaces. *IEEE Transactions on Information Theory*, 56(11):5839–5846, 2010.
- Colin McDiarmid. On the method of bounded differences. Surveys in combinatorics, 141(1):148–188, 1989.
- Nishant A. Mehta and Alexander G. Gray. On the sample complexity of predictive sparse coding. arXiv preprint arXiv:1202.4050v1, 2012.
- Ron Meir and Tong Zhang. Generalization error bounds for bayesian mixture algorithms. *Journal of Machine Learning Research*, 4:839–860, 2003.
- Shahar Mendelson and Petra Philips. On the importance of small coordinate projections. *Journal of Machine Learning Research*, 5:219–238, 2004.
- Michael R. Osborne, Brett Presnell, and Berwin A. Turlach. On the lasso and its dual. *Journal of Computational and Graphical statistics*, pages 319–337, 2000.
- John Shawe-Taylor, Peter L. Bartlett, Robert C. Williamson, and Martin Anthony. Structural risk minimization over data-dependent hierarchies. *IEEE Transactions on Information Theory*, 44 (5):1926–1940, 1998.
- David Slepian. The one-sided barrier problem for Gaussian noise. *Bell System Technical Journal*, 41(2):463–501, 1962.
- Ingo Steinwart and Andreas Christmann. Support vector machines. Springer, 2008.
- Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996.
- Daniel Vainsencher, Shie Mannor, and Alfred M. Bruckstein. The sample complexity of dictionary learning. *Journal of Machine Learning Research*, 12:3259–3281, 2011.
- Vladimir N. Vapnik and Alexey Ya. Chervonenkis. Uniform convergence of frequencies of occurence of events to their probabilities. In *Dokl. Akad. Nauk SSSR*, volume 181, pages 915–918, 1968.
- Mathukumalli Vidyasagar. Learning and Generalization with Applications to Neural Networks. Springer, London, 2002.

Kai Yu, Tong Zhang, and Yihong Gong. Nonlinear learning using local coordinate coding. In Yoshua Bengio, Dale Schuurmans, John Lafferty, Christopher K. I. Williams, and Aron Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 2223–2231. MIT Press, 2009.

Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):301–320, 2005.

G Glossary

Notation	Description	Page
		List
D_j	j^{th} atom (column) of dictionary (matrix) D	1
[n]	$ \{1,2,\ldots,n\} $	4
\mathcal{D}	The space of dictionaries $(B_{\mathbb{R}^d})^k$	3
\mathcal{D}_{μ}	$\{D \in \mathcal{D} : \mu_s(D) \ge \mu\}$	13
$egin{array}{c} \mathcal{D}_{\mu} \ \mathcal{F} \end{array}$	$ \{ f_{D,w} := x \mapsto \langle w, \varphi_D(x) \rangle : D \in \mathcal{D}, w \in \mathcal{W} \}.$	3
\mathcal{F}_{μ}	$\mid \{f = (D, w) \in \mathcal{F} : D \in \mathcal{D}_{\mu}\}$	13
$ \begin{array}{l} \mathcal{F}_{\mu} \\ \mathcal{F}_{\mu^*,\eta}(\mathbf{x}) \\ \mathcal{F}_{\mu^*} \\ \mathcal{F}_{\mu^*}(\mathbf{x}) \\ \Pi \end{array} $	$\left\{ f \in \mathcal{F}_{\mu^*} : \exists \tilde{\mathbf{x}} \subseteq_{\eta} \mathbf{x} \text{ s-sparse}(\varphi_D(\tilde{\mathbf{x}})) \text{ and } \left[\operatorname{margin}_s(D, \tilde{\mathbf{x}}) > \tau \right] \right\}$	20
\mathcal{F}_{μ^*}	$\{f = (D, w) \in \mathcal{F} : (\mu_s(D) \ge \mu_s^*) \text{ and } (\mu_{2s}(D) \ge \mu_{2s}^*) \}$	19
$\mathcal{F}_{\mu^*}(\mathbf{x})$	$\{f \in \mathcal{F}_{\mu^*} : s\text{-sparse}(\varphi_D(\mathbf{x})) \text{ and } [\text{margin}_s(D, \mathbf{x}) > \tau]\}$	20
	Marginal probability measure over input space $B_{\mathbb{R}^d}$	3
P	Joint probability measure over $B_{\mathbb{R}^d} \times \mathcal{Y}$	3
$Pl(\cdot, f)$	$\mid E_{(x,y)} l(y(f(x)))$	7
Pf	$\mid E_{(x,y)\sim P}f(x)$	7
$P_{\mathbf{z}}l(\cdot,f)$	$\begin{bmatrix} E_{(x,y)\sim P}f(x) \\ \sum_{i=1}^m l(y_i,f(x_i)) \\ \sum_{i=1}^m f(x_i) \end{bmatrix}$	7
$P_{\mathbf{z}}f$	$\sum_{i=1}^{m} f(x_i)$	7
\mathcal{W}	The space of linear hypotheses, equal to $rB_{\mathbb{R}^d}$	3
\mathcal{Y}	Space of labels or targets	3
$lpha B_{\mathbb{R}^d}$	The ball in \mathbb{R}^d of radius α	3
$\hat{f}_{\mathbf{z}}$	Hypothesis returned by learner from \mathbf{z}	7
γ_i	standard normal random variable	10
$\mathcal{G}_{m \mathbf{x}}(\mathcal{F})$	conditional Gaussian average of \mathcal{F}	11
$\tilde{\mathbf{x}} \subseteq_{\eta} \mathbf{x}$	$\tilde{\mathbf{x}}$ is a subset of \mathbf{x} with at most η elements of \mathbf{x} removed	12
\mathbf{x}''	unlabeled m-sample	8
\mathbf{z}'	second labeled <i>m</i> -sample (ghost sample)	10
Z	labeled m-sample of training data	3
$\mu_s(D)$	s-incoherence: minimum $(\varsigma_s(D))^2$ among s-atom subdictionaries of D	4
$ S _{2,s}$	The s-restricted 2-norm: $\sup_{t \in \mathbb{R}^n: t =1, \operatorname{supp}(t) \leq s} St _2$	21
$\mathcal{R}_{m \mathbf{x}}(\mathcal{F})$	conditional Rademacher average of \mathcal{F}	11
σ_i	Rademacher random variable	10
$\operatorname{margin}_s(D, \mathbf{x})$	$\min_{x_i \in \mathbf{x}} \operatorname{margin}_s(D, x_i)$	5
$\operatorname{margin}_s(D, x)$	$\left \max_{\substack{\mathcal{I} \subseteq [k] \\ \mathcal{I} = k - s}} \min_{j \in \mathcal{I}} \left\{ \lambda - \left \langle D_j, x_i - D\varphi_D(x_i) \rangle \right \right\} \right $	5
supp(t)	$ \mid \{i \in [k] : t_i \neq 0\} $	4
$arphi_D$	$\varphi_D(x) := \arg\min_z \ x - Dz\ _2^2 + \lambda \ z\ _1$	3
$\varsigma_s(A)$	the s^{th} singular value of A	4
s -sparse $(\varphi_D(\mathbf{x}))$	for all $x_i \in \mathbf{x}$, $\ \varphi_D(x_i)\ _0 \le s$	4