**DENOISING BY SPARSE APPROXIMATION:
ERROR BOUNDS BASED ON
RATE-DISTORTION THEORY**


by

Alyson K. Fletcher, Sundeep Rangan, Vivek K. Goyal
and Kannan Ramchandran

# DENOISING BY SPARSE APPROXIMATION:
# ERROR BOUNDS BASED ON
# RATE-DISTORTION THEORY

by

Alyson K. Fletcher, Sundeep Rangan, Vivek K. Goyal
and Kannan Ramchandran

Memorandum No. UCB/ERL M05/5

23 September 2004

## ELECTRONICS RESEARCH LABORATORY

College of Engineering
University of California, Berkeley
94720

# Denoising by Sparse Approximation:
# Error Bounds Based on Rate–Distortion Theory*

Alyson K. Fletcher,[†] Sundeep Rangan,[‡]
Vivek K Goyal,[§] and Kannan Ramchandran

## Abstract

If a signal $x$ is known to have a sparse representation with respect to a frame, the signal can be estimated from a noise-corrupted observation $y$ by finding the best sparse approximation to $y$. The ability to remove noise in this manner depends on the frame being designed to efficiently represent the signal while it *inefficiently* represents the noise. This paper gives bounds on the expected squared error of this denoising scheme. The main challenge in computing the expected squared error is that the noise realization affects the choice of subspace in the sparse approximation of $y$.

One error bound depends on the expected fraction of energy of a white Gaussian signal in a signal-dependent choice of subspaces. A bound on this expected fraction, derived using rate–distortion theory, may be of independent interest. Furthermore, for certain randomly generated frames, a simple expression is obtained for the probability that the estimate lies in the correct subspace.

**Keywords:** beta function, denoising, dictionary, estimation, F-distribution, Gamma function, Gaussian signals, Grassmannian frames, information theory, maximum likelihood, nonlinear approximation, Occam filters, rate–distortion, sparse approximation, stable signal recovery, subspace fitting

---

[†]A. K. Fletcher (alyson@eecs.berkeley.edu) and K. Ramchandran (kannanr@eecs.berkeley.edu) are with the Dept. of Electrical Engineering and Computer Sciences, University of California, Berkeley, CA 94720 USA.

[‡]S. Rangan (s.rangan@flarion.com) is with Flarion Technologies, Bedminster, NJ 07921 USA.

[§]V. K. Goyal (vgoyal@mit.edu) is with the Dept. of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA 02139 USA.

# 1  Introduction

## 1.1  Denoising by Sparse Approximation with a Frame

Consider the problem of estimating an unknown signal $x \in \mathbb{R}^N$ from the noisy observation $y = x + d$ where $d \in \mathbb{R}^N$ has the i.i.d. Gaussian $\mathcal{N}(0, \sigma^2 I_N)$ distribution. If $x$ is known to lie in a given $K$-dimensional subspace of $\mathbb{R}^N$, the situation can immediately be improved by projecting $y$ to the given subspace; since the noise distribution is spherically symmetric and the projection is independent of the noise, only $K/N$ fraction of the original noise is then left. Further information about the distribution of $x$ could be exploited to remove even more noise.

In this paper we consider the estimation of $x$ from $y$ with a weaker, more permissive, signal model. Rather than knowing a $K$-dimensional subspace that contains $x$, we are given a *set of $K$-dimensional subspaces* such that $x$ is contained in their union. Specifically, let $\Phi = \{\varphi_i\}_{i=1}^M \subset \mathbb{R}^N$, $M \geq N$, be a frame. Our model is that

$$x = \sum_{i=1}^M \alpha_i \varphi_i$$

for some coefficient vector $\alpha \in \mathbb{R}^M$ with at most $K$ nonzero elements. An equivalent statement is that $x$ lies in one of the $J = \binom{M}{K}$ subspaces obtained by selecting a $K$-element subset of $\Phi$. When $K \ll N$, it is said that $x$ is *represented sparsely* in the frame $\Phi$.

With the addition of the noise $d$, the observed vector $y$ will (almost surely) not be represented sparsely. Intuitively, the point from one of the $J$ subspaces under consideration that is closest to $y$ is a good estimate for $x$. Formally, because the probability density function of $d$ is a strictly decreasing function of $\|d\|$, this is the maximum likelihood estimate of $x$ given $y$. We will write

$$\hat{x} = \operatorname*{argmin}_{\left\{ x \;\mid\; x = \sum_{i=1}^M \alpha_i \varphi_i \text{ with at most } K \text{ nonzero } \alpha_i \text{s} \right\}} \|y - x\|_2 \tag{1}$$

for this estimate and call it the optimal $K$-term approximation of $y$. Henceforth we omit the subscript 2 that indicates the use of the Euclidean norm.

The main results of this paper are bounds on the per-component mean-squared estimation error $\frac{1}{N}\mathbf{E}\left[\|x - \hat{x}\|^2\right]$ for denoising via sparse approximation. (The expectation is over the noise $d$ only, as we do not assume a probabilistic model for $x$.) These bounds depend on $(N, M, K)$ but avoid further dependence on the frame $\Phi$ (such as the coherence of $\Phi$); some results hold for all $\Phi$ and others are for randomly generated $\Phi$. To the best of our knowledge, the results are novel for

2

(a) being an average-case (rather than worst-case) analysis;

(b) having dependence on frame size rather than more fine-grained properties of the frame; and

(c) using source coding gedanken experiments in this context.

Preliminary results were first presented in [16].

## 1.2 Connections to Approximation and to Practice

A likely situation in practice is that the underlying true signal $x$ has a good $K$-term *approximation* rather than an exact $K$-term *representation*. At very least, this is the goal in designing the frame $\Phi$ for a signal class of interest. It is then still reasonable to compute (1) to estimate $x$ from $y$, but there are trade-offs in the selections of $K$ and $M$.

Let $f_{M,K}$ denote the squared Euclidean approximation error of the optimal $K$-term approximation, where the subscript $M$ emphasizes the frame size. It is obvious that $f_{M,K}$ decreases with increasing $K$, and with suitably designed frames it also decreases with increasing $M$. One concern of approximation theory is to study the decay of $f_{M,K}$ precisely. (For this we should consider $N$ very large or infinite.) For piecewise smooth signals, for example, wavelet frames give exponential decay with $K$ [5, 10, 13].

When one uses sparse approximation to denoise, the performance depends on both the ability to approximate $x$ and the ability to reject the noise. Approximation is improved by increasing $M$ and $K$, but noise rejection is diminished. The dependence on $K$ is clear, as the fraction of the original noise that remains on average is at least $K/N$. For the dependence on $M$, note that increasing $M$ increases the number of subspaces and thus increases the chance that the selected subspace is not the best one for approximating $x$. When $M$ is very large, there is some subspace very close to $y$ and thus $\hat{x} \approx y$.

Fortunately, there are many classes of signals for which $M$ need not grow too quickly as a function of $N$ to get good sparse approximations. The design of frames $\Phi$ for this problem is essentially the same as designing dictionaries for matching pursuit [23]. Examples of audio dictionaries with good computational properties were given by Goodwin [18]. See also Moschetti *et al.* [24] for video compression and Engan *et al.* [15] for an iterative design procedure.

One motivation for this work was to give guidance for the selection of $M$. This requires the combination of approximation results (e.g., bounds on $f_{M,K}$) with results such as ours. The results presented here do not address approximation quality.

## 1.3 Related Work

Computing optimal $K$-term approximations is generally a difficult problem. Given $\epsilon \in \mathbb{R}^+$ and $K \in \mathbb{Z}^+$, to determine if there exists a $K$-term approximation such that

$\|x - \hat{x}\| \leq \epsilon$ is an NP-complete problem [9, 26]. Similarly, given $K \in \mathbb{Z}^{+}$, to find an $\alpha$ that minimizes $\|x - \hat{x}\|$ among all vectors with not more than $K$ nonzero entries is NP-hard [9].

The difficulty of optimal sparse approximation has prompted study of heuristics. A greedy heuristic that is standard for finding sparse approximate solutions to linear equations [17] has been known as *matching pursuit* in the signal processing literature since the work of Mallat and Zhang [23]. One of Mallat and Zhang's initial applications was the separation of signal from noise. Chen, Donoho and Saunders [3, 4] proposed a convex relaxation of the approximation problem (1) called *basis pursuit*.

Both matching pursuit and basis pursuit have been the subject of many analytical and numerical investigations. Recently, Donoho, Elad, Temlyakov, and Tropp have determined conditions on the frame such that matching pursuit and basis pursuit find optimal approximations [11, 12, 31, 32]. An overlapping body of recent literature establishes that optimal approximations have stability properties, depending on a coherence measure of the frame, that make the positions of the nonzero entries in $\alpha$ insensitive to additive noise up to some bound in magnitude [11, 14, 20]. The results of this paper look beyond the noise magnitude threshold at which the footprint of the optimal $\alpha$ is unchanged.

Denoising by finding a sparse approximation is similar to the concept of denoising by compression popularized by Saito [28] and Natarajan [25]. More recent works in this area include those by Krim *et al.* [21], Chang *et al.* [2] and Liu and Moulin [22]. All of these works use bases rather than frames. To put the present work into a similar framework would require a "rate" penalty for redundancy. Instead, the only penalty for redundancy comes from choosing a subspace that does not contain the true signal ("overfitting" or "fitting the noise"). The literature on compression with frames notably includes [1, 19, 27].

This paper uses quantization and rate–distortion theory only as a proof technique; there are no encoding rates because the problem is purely one of estimation. However, the "negative" results on representing white Gaussian signals with frames presented here should be contrasted with the "positive" encoding results of Goyal *et al.* [19]. The positive results are limited to low rates (and hence signal-to-noise ratios that are usually uninteresting). A natural extension of this work is to derive negative results for encoding. This would support the assertion that frames in compression are useful not universally, but only when they can be designed to yield very good sparseness for the signal class of interest.

## 1.4 Preview of Results

To motivate the paper, we present a set of numerical results from Monte Carlo simulations that qualitatively reflect our main results. The sizes of $N$, $M$, and $K$ are small because of the high complexity of optimal approximation and because a large number
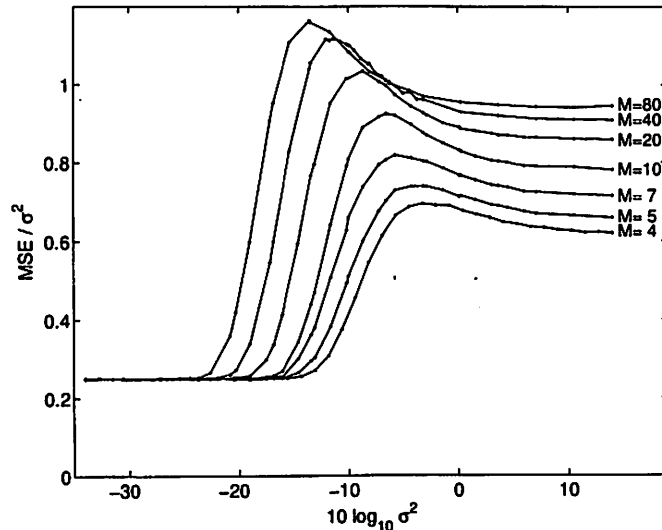
Figure 1: Performance of denoising by sparse approximation when the true signal $x \in \mathbb{R}^4$ has an exact 1-term representation with respect to a frame that is an optimal $M$-element Grassmannian packing.

of independent trials is needed to get adequate precision. Each data point shown is the average of 100 000 trials.

Consider a true signal $x \in \mathbb{R}^4$ ($N = 4$) that has an exact 1-term representation ($K = 1$) with respect to $M$-element frame $\Phi$. We observe $y = x+d$ with $d \sim \mathcal{N}(0, \sigma^2 I_4)$ and compute estimate $\hat{x}$ from (1). The signal is generated with unit norm so that the signal-to-noise ratio (SNR) is $1/\sigma^2$ or $-10\log_{10}\sigma^2$ dB.

To have tunable $M$, we used frames that are $M$ maximally separated unit vectors in $\mathbb{R}^N$, where separation is measured by the minimum pairwise angle among the vectors and their negations. These are cases of Grassmannian packings [6, 30] in the simplest case of packing one-dimensional subspaces (lines). We used packings tabulated by Sloane with Hardin, Smith and others [29].

Throughout we use the following definition for mean-squared error:

$$\text{MSE} = \frac{1}{N}\mathbf{E}\left[\|x - \hat{x}\|^2\right].$$

Fig. 1 shows the MSE as a function of $\sigma$ for several values of $M$. Note that for visual clarity, MSE$/\sigma^2$ is plotted. For small values of $\sigma$, the MSE is $(1/4)\sigma^2$. This is an example of the general statement that

$$\text{MSE} = \frac{K}{N}\sigma^2 \qquad \text{for small } \sigma,$$
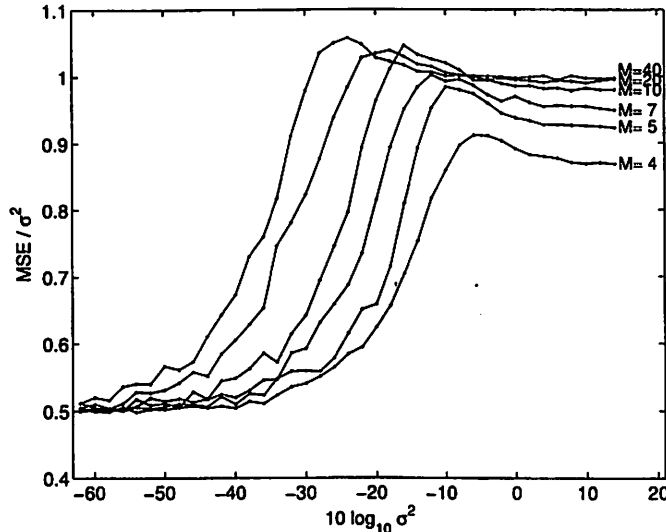
5

Figure 2: Performance of denoising by sparse approximation when the true signal $x \in \mathbb{R}^4$ has an exact 2-term representation with respect to a frame that is an optimal $M$-element Grassmannian packing.

as described in detail in Section 2. For large values of $\sigma$, the scaled MSE approaches a constant value:

$$\lim_{\sigma \to \infty} \frac{\text{MSE}}{\sigma^2} = g_M,$$

where $g_M$ is a slowly increasing function of $M$ and $\lim_{M \to \infty} g_M = 1$. The characterization of the dependence on $M$ is the main contribution of Section 3. Finally, the transition between low and high SNR is studied in Section 4. The same properties are illustrated for $K = 2$ in Fig. 2.

# 2    Preliminary Computations

The set of vectors comprising the dictionary is denoted $\Phi$. In the following analysis, it is possible for the dictionary to be undercomplete, i.e. $M < N$. The model for the signal $x$ is that it is known to be some linear combination of $K$ of the $M$ dictionary vectors. The $K$ vectors to select and the $K$ coefficients of the linear combination must be estimated from $y$. Furthermore, it is assumed that $K \leq N$; otherwise there is no sparseness at all.

One could use any of several techniques to form the estimate $\hat{x}$ from $y$. Without introducing a probability distribution for $x$, the maximum likelihood estimate over the

noise $d$ is considered. This yields

$$\hat{x}_{ML}(y) \;=\; \operatorname*{argmin}_{x \in X} \|y - x\|, \tag{2}$$

where $X$ is the set of all $x$ spanned by $K$ of the $M$ frame vectors in $U$. In this scenario, the error of the ML estimate provides a *lower* bound on the average error achievable by any other estimator. Many methods used in practice attempt to emulate the ML estimator.

The view of the ML estimate as a projection is used in the analysis. Given a frame $U$, let $\mathcal{P}_K$ be the set of all projections onto spaces spanned by $K$ of the $M$ frame vectors in $U$. Then $\mathcal{P}_K$ has $\binom{M}{K}$ projections, and the ML estimate is given by

$$\hat{x}_{ML} \;=\; \hat{P}y, \tag{3}$$

where

$$\hat{P} \;=\; \operatorname*{argmax}_{P \in \mathcal{P}_K} \|Py\|.$$

We now analyze the error of the ML estimator. The error usually depends on the true signal $x$. Define the conditional MSE

$$e(x) \;=\; \frac{1}{N}\mathbf{E}\left(\|x - \hat{x}_{ML}\|^2 \mid x\right).$$

When $M = K$, there is no subspace selection: there is just one approximation subspace. In this case, the ML estimate $\hat{x}_{ML}$ is the projection of the noisy signal onto a fixed $K$-dimensional subspace, which is equivalent to standard least squares (LS) denoising. The estimation error does not then depend on $x$. The error reduces by a factor $K/N$,

$$e(x) \;=\; \frac{K}{N}\sigma^2 \qquad \forall x. \tag{4}$$

At high SNR (small $\sigma^2$), the selection of the $K$-dimensional subspace is unperturbed by $d$. The error expression (4) thus holds at this extreme as well.

# 3 Rate-Distortion Analysis

## 3.1 Sparse Approximation of a Gaussian Source

Before directly addressing the denoising performance of sparse approximation, we give a new approximation result for white Gaussian signals. This result is a *lower bound* on the error in sparse approximation of a white Gaussian signal; it serves as the basis for *upper bounds* on the MSE for denoising when the SNR is low.

**Theorem 1** *Let $\Phi$ be an M-element frame and let $d \in \mathbb{R}^N$ be a zero-mean, Gaussian vector with variance $\sigma^2 I_N$. For any $\hat{d}$ that can be represented as a linear combination of K vectors from $\Phi$,*

$$\frac{1}{N} E\left[\|d - \hat{d}\|^2\right] \geq \sigma^2 J^{-1/r} \left( \left(\frac{K}{N}\right)^{K/(N-K)} - \left(\frac{K}{N}\right)^{N/(N-K)} \right), \qquad (5)$$

*where*

$$J = \binom{M}{K} \qquad and \qquad r = \frac{N-K}{2}. \qquad (6)$$

*For large N and small K/N, the bound reduces to the simple expression*

$$\frac{1}{N} E\left[\|d - \hat{d}\|^2\right] \geq \sigma^2 J^{-1/r}(1 - K/N). \qquad (7)$$

**Proof:** See Appendix A.

Theorem 1 shows that for any $\Phi$, there is an approximation error lower bound that depends only on the frame size $M$, the dimension of the signal $N$, and the dimension of the signal model $K$. Henceforth we use the notation

$$D_{\min} = \min\left(\frac{1}{N} E\left[\|d - \hat{d}\|^2\right]\right) \qquad (8)$$

where the minimization is over all functions that map $d$ to a $K$-sparse vector $\hat{d}$.

Observe that when $M = K$, the bound in (7) reduces to

$$D_{\min} \geq \sigma^2(1 - K/N),$$

which is tight: When $M = K$ there is only one subspace and the frame will capture a fraction $K/N$ of the source signal power.

At the other extreme, as $M \to \infty$, $J \to \infty$, and $D_{\min} \to 0$. This limit is also expected: as $M \to \infty$, the frame can eventually capture all the power of the Gaussian signal.

However, when $K/N$ is small, the error $D_{\min}$ reduces to zero very slowly as $M$ increases. To see this, define a sparsity measure $\alpha = K/N$ and a redundancy factor $\rho = M/N$. Then for large $N$, the bound (7) reduces to

$$D_{\min} \geq \sigma^2 \left(\frac{\alpha}{\rho}\right)^\alpha \left(1 - \frac{\alpha}{\rho}\right)^{\rho - \alpha} (1 - \alpha). \qquad (9)$$

For a large redundancy $\rho \to \infty$, $D_{\min} = O(\rho^{-\alpha})$. Consequently, $D_{\min} \to 0$ slowly when $\alpha$ is small. In other words, when the sparsity $K/N$ is small, an exponentially large number of frame vectors are needed to capture the signal. In this way, we can say that a white Gaussian signal is not sparsely approximated well by any frame.

8

## 3.2 Comparison to Simulated Random Frames with $K = 1$

The bound in Theorem 1 is not, in general, tight: the inequality in the theorem provides only a lower bound on the error. The actual error will depend on the specific frame. However, we expect the bound to be tight for frames designed specifically to represent the Gaussian signal, and the qualitative value of the bound is supported by the Monte Carlo simulations described in this section.

To empirically evaluate the tightness of the bound, we compare it to the actual error obtained using random frames. We will address random frames in greater detail in Section 4. For a random frame, we first fix some value for $N$. Then, for various values of $M$, a random frame is obtained by generating a random $M \times N$ matrix $A$ and then finding $U$, the $M \times N$ matrix whose columns are the $M$ dominant orthonormal eigenvectors of $AA'$. The matrix $U$ can be found by a singular value decomposition of $A$. The $M$ frame vectors are then given by the $M$ rows of $U$.

With each random frame, we find the best $(K = 1)$-dimensional approximation, $\hat{d}$, to a random $N$-dimensional vector, $d$, whose components are i.i.d. Gaussian with zero-mean and variance $\sigma^2 = 1$. Using 100 random frames, we estimate the average value $\mathbf{E}\left[\|\hat{d}\|^2\right]$. Fig. 3 plots the normalized error

$$-10 \log_{10} \frac{\mathbf{E}\|d - \hat{d}\|^2}{N\sigma^2} \qquad (10)$$

as a function of the frame sizes $M$ for signal lengths $N = 10$ and $N = 100$. Also plotted in Fig. 3 is the theoretical lower bound from Theorem 1.

## 3.3 Bound on MSE

We now return to our original denoising problem defined in Section 2. We wish to bound the estimation error $e(x)$ for a given signal $x$ and frame $\Phi$.

Theorem 1 provides a lower bound on the ability of a frame to capture white noise with a $K$-term approximation. Intuitively, if the true signal $x$ is perfectly modeled with such an approximation, frame-based denoising should capture the signal and remove the noise. We now quantitatively justify this assertion.

To state the result, we need some more notation. As discussed earlier, there are $J = \binom{M}{K}$ subspaces spanned by $K$ of the $M$ frame vectors in $\Phi$. Index the subspaces by $j = 1, \ldots, J$ and let $T$ be the index of the subspace closest to the noisy signal $y$. Thus, the ML estimator $\hat{x}_{ML}$ is the projection of $y$ onto the subspace of index $T$. For a fixed true signal, $x$, the subspace selection variable, $T$, is a discrete random variable depending on the noisy signal $y$. Let $H(T)$ denote its entropy in bits.

**Theorem 2** *Consider the denoising problem described in Section 2 with a fixed true signal $x$. Assume the true signal can be represented by $K$ of the $M$ frame vectors in $\Phi$.*
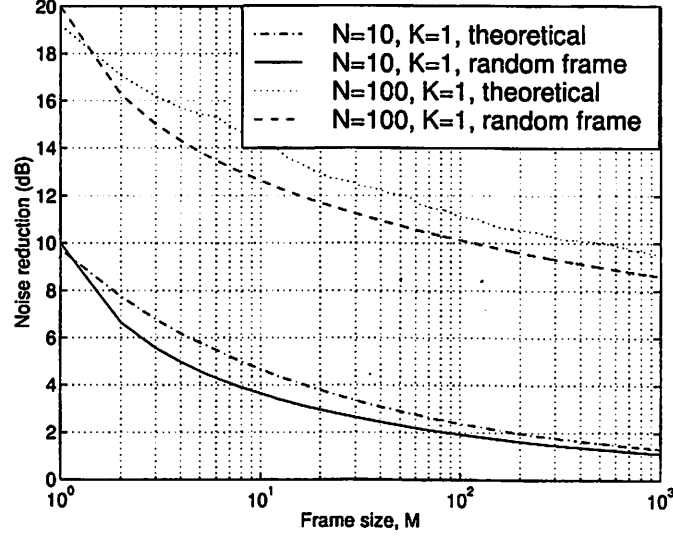
9

Figure 3: Comparison of theoretical noise-fit lower bound to actual noise-fit error with a random frame.

Let $H(T)$ be the entropy of the subspace selection variable. Assume the ML estimator is unbiased. That is,

$$E[\hat{x}_{ML} \mid x] = x.$$

Then, the MSE is bounded above by

$$e(x) \leq \sigma^2 \left( 1 - 2^{-H(T)/r} \left( 1 - \frac{K}{N} \right) \right). \tag{11}$$

**Proof:** See Appendix B.

Theorem 2 requires that the ML estimate is unbiased. Unfortunately, this does not generally hold. The frame-based estimator may project to the wrong subspace and lose signal energy. These incorrect projections will result in an estimator whose average value is in general slightly smaller than the true signal. The theorem therefore represents an approximation of the actual situation and is useful when combined with an estimate of the bias error with sparse approximation.

This theorem requires the entropy $H(T)$. Evaluation of $H(T)$ may be problematic; however, we can bound the entropy in two simple cases.

First, since $T$ takes on values in $\{1, 2, \ldots, J\}$, $H(T)$ is upper bounded by $\log_2 J$. This gives the bound

$$e(x) \leq \sigma^2 \left( 1 - J^{-1/r} \left( 1 - \frac{K}{N} \right) \right). \tag{12}$$

10

As a second bound, suppose we know $p_\text{corr}$, the probability that the estimator selects the correct subspace. Then, the maximum entropy for $T$ occurs when $T$ has probability $p_\text{corr}$ on the correct subspace and uniform probabilities $(1 - p_\text{corr})/(J - 1)$ on the remaining $J - 1$ subspaces. Thus, the entropy is upper bounded by,

$$
\begin{aligned}
H(T) &\leq -p_\text{corr} \log_2 p_\text{corr} - (J - 1)\frac{(1 - p_\text{corr})}{J - 1} \log_2\left(\frac{(1 - p_\text{corr})}{J - 1}\right) \\
&= H(p_\text{corr}) + (1 - p_\text{corr})\log_2(J - 1),
\end{aligned}
$$

where $H(p)$ is the binary entropy

$$
H(p) = -p \log_2 p - (1 - p)\log_2(1 - p).
$$

Substituting this in (11) gives

$$
e(x) \leq \sigma^2 \left( 1 - 2^{-H(p_\text{corr})/r}(J - 1)^{-(1 - p_\text{corr})/r}\left(1 - \frac{K}{N}\right)\right). \tag{13}
$$

The bound shows that as $p_\text{corr} \to 1$, $e(x) \to \sigma^2 K/N$, which is the result when denoising to a single subspace.

# 4  Large Random Frame Analysis

## 4.1  Random Frames

The expression for the MSE in the previous section is, in general, difficult to evaluate for an arbitrary frame. In order to obtain more concrete results, in this section, we consider large *random* frames where the analysis is significantly simpler.

Specifically, we assume that the frame $\Phi$ consists of a large number of independent random vectors $\varphi_i$ uniformly distributed on the unit sphere. We will compute the MSE averaged over the random frame.

As before, we will assume that the true signal $x$ is represented exactly by a linear combination of $K$ of the $M$ frame vectors. We will take $\|x\| = N$ so that the SNR is

$$
\gamma = \frac{\|x\|^2}{\mathbf{E}\left[\|d\|^2\right]} = \frac{N}{N\sigma^2} = \frac{1}{\sigma^2}.
$$

## 4.2  Probability of Correct Subspace Selection

We first compute the probability of selecting the correct subspace using the ML estimator. To state the main result, we need the following function. For $a, b, r \geq 0$ define

$$
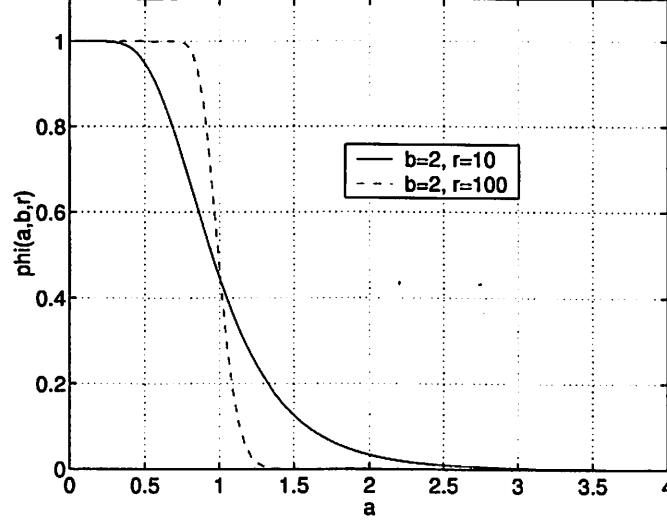\phi(a, b, r) = \int_0^\infty p_r(u)\exp\left(-\left(\frac{abu}{1 + au}\right)^r\right)du \tag{14}
$$

Figure 4: Function $\phi(a, b, r)$ for $b = 2$ and various values of $a$.

where

$$p_r(u) = r^r \Gamma(r) u^{r-1} e^{-ru} \tag{15}$$

and $\Gamma(r)$ is the Gamma function [36]. The function $\phi$ can be computed numerically as a single line integral. The following lemma describes the function qualitatively.

**Lemma 1** *For the function $\phi(a, b, r)$ defined above:*

(a) *For all $a$,$b$ and $r$, $\phi(a, b, r) \in [0, 1]$.*

(b) *For all $r$, $\phi(a, b, r)$ monotonically decreases in $a$ and $b$ with $\phi(a, 0, r) = \phi(0, b, r) = 1$, and*

$$\lim_{b \to \infty} \phi(a, b, r) = 0, \quad \lim_{a \to \infty} \phi(a, b, r) = 0.$$

(c) *In the limit as $r \to \infty$, $\phi(a, b, r)$ is a step function. That is, for all $a$ and $b$,*

$$\lim_{r \to \infty} \phi(a, b, r) = \begin{cases} 1, & \text{if } ab < 1 + a; \\ 0, & \text{if } ab > 1 + a. \end{cases}$$

To illustrate the function, Fig. 4 shows $\phi(a, b, r)$ as a function of $a$ for $b = 2$ and $r = 10$ and 100. As can be seen, $\phi(a, b, r)$ monotonically decreases in $a$. Also, the transition from 1 to 0 becomes sharper at the higher value of $r$.

12

Now consider the frame problem described earlier. Let

$$r = \frac{N-K}{2}, \quad a = \frac{N-K}{N}, \quad C = \left[\frac{J}{r\beta(r, K/2)}\right]^{1/r}, \tag{16}$$

where $\beta(p, q)$ is the beta function [33]. We now state the main result of this section.

**Theorem 3** *For sparse approximation with respect to a random frame, if $M \gg K$ and $N \gg K$, the probability of selecting the correct subspace is given by*

$$p_{corr} = \phi(a\sigma^2, C, r).$$

**Proof:** See Appendix C.

The result shows that, for large random frames, the probability of selecting the correct subspace can be computed with a single line integral. Also, the asymptotic expression in Lemma 1(c) shows

$$\lim_{r\to\infty} p_{corr} = \begin{cases} 1, & \text{if } aC\sigma^2 < 1 + a\sigma^2; \\ 0, & \text{if } aC\sigma^2 > 1 + a\sigma^2. \end{cases}$$

If we define the *critical* SNR,

$$\gamma_{crit} = a(C-1),$$

the asymptotic expression can be rewritten as

$$\lim_{r\to\infty} p_{corr} = \begin{cases} 1, & \text{if } \gamma > \gamma_{crit}; \\ 0, & \text{if } \gamma < \gamma_{crit}, \end{cases}$$

where $\gamma = 1/\sigma^2$ is the SNR. Thus, for large $r$, there is a critical SNR with a simple expression where the probability of correct subspace selection transitions from 0 to 1.

## 4.3 MSE Bound

The previous subsection estimates the probability $p_{corr}$ that the ML estimator selects the correct subspace. However, in general, we are interested in the MSE of the estimator. One simple method to estimate the MSE is to substitute the estimate for $p_{corr}$ into (13). However, there are two difficulties with this approach.

First, as stated earlier, an assumption in Theorem 2 is that the ML estimator is unbiased. Actually, the ML estimator (1) is only approximately unbiased.

Secondly, Theorem 2 technically holds for a specific frame with a specific selection probability $p_{corr}$. To obtain the correct expected MSE, for each frame, we should compute its $p_{corr}$, evaluate the MSE in Theorem 2, and then average over all the frames. However, $p_{corr}$ in Theorem 3 is the selection probability already averaged over the frame. Substituting the average selection probability $p_{corr}$ into the MSE expression is not the same as averaging over the MSEs.

Nevertheless, we will show in the numerical experiment in the next subsection that the simple substitution of $p_{corr}$ into (13) appears to give a reasonable estimate for the MSE. However, the theoretical basis for this requires further investigation.
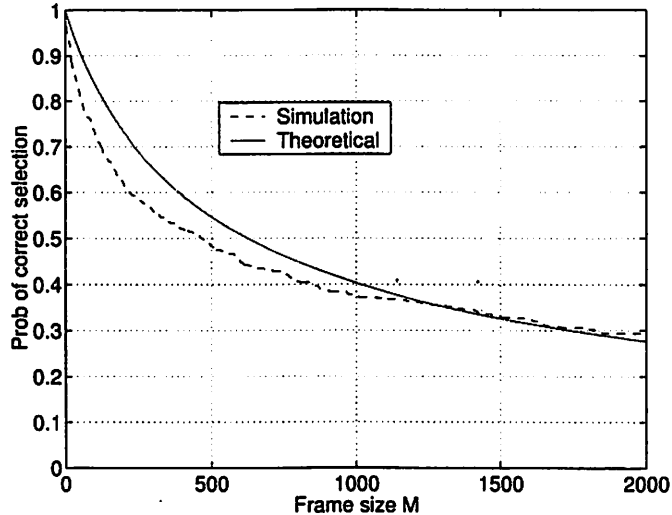
13

Figure 5: Average probability of selecting the correct subspace with $N = 10$, $K = 1$ and random frames of various sizes $M$. The measured selection probability is compared against the estimate in Theorem 3.

## 4.4 Numerical Simulation

As a simple numerical experiment, we evaluated the estimates of the previous subsections with signal dimension $N = 10$. In the space $\mathbb{R}^N$, we generated various random frames consisting of $M$ vectors uniformly distributed on the unit sphere. The frame size $M$ was varied from 1 to 1000. The true signal $x$ is taken to be one of $K = 1$ of the $M$ vectors. Noise $d$ was added at an SNR of 3 dB, and the maximum likelihood estimator $\hat{x}_{ML}$ was computed as described in Section 2.

Fig. 5 shows the average probability that the estimator selects the correct subspace. For each frame size $M$, the probability is averaged over 500 random frames. Also plotted is the theoretical error probability from Theorem 3. The figures shows that as $M$ increases the theoretical value matches the simulated value closely.

Fig. 6 shows the average noise reduction for the same problem. Noise reduction is defined as

$$-10 \log_{10} \left[ \frac{\mathbf{E} \|x - \hat{x}_{ML}\|^2}{N\sigma^2} \right]$$

which represents the error normalized by the original noise value. Also plotted in Fig. 6 is an estimate of the noise reduction based on substituting the subspace selection probability from Theorem 3 into (13). As stated in Section 4.3, this substitution is not theoretically correct due to variation of the selection probability. Also, Theorem

14

Figure 6: Average noise reduction with $N = 10$, $K = 1$ and random frames of various sizes $M$. The measured selection noise reduction is compared against a theoretical value.

2 does not quite apply since the estimate $\hat{x}_{ML}$ is not unbiased. To account for this fact, we measured the bias error in the simulation and added it to the MSE from Theorem 2. With this modification, the theoretical noise reduction matches the measured value within 0.5 dB. However, properly accounting for the bias error needs further investigation.

# A  Proof of Theorem 1

There are at most $J$ possible projections in $\mathcal{P}_K$. Denote the projections by $\{P_j\}_{j=1}^J$, and let $\{R_j\}_{j=1}^J$ be the corresponding range spaces.

The minimization in (8) is given by

$$D_{\min} = \frac{1}{N}\mathbf{E}\|d - \hat{d}\|^2 \qquad (17)$$

where

$$\hat{d} = P_T d, \quad \text{where } T = \operatorname*{argmax}_j \|P_j d\|. \qquad (18)$$

Here, $T$ is the index of the subspace with the most energy of $d$. The minimization is given by the projection of $d$ onto that subspace.

Now, fix any positive real number $n$. For each $j \in \{1, 2, \ldots, J\}$, let $Q_j$ be an optimal $n$-bit quantizer of $P_T d$ conditional on $T = j$. Define the quantizer $Q$ by

$$Q(d) = Q_T(P_T(d)).$$

That is, $Q$ projects $d$ to the $K$-dimensional space with the greatest energy and then quantizes the projection of $d$ within that space.

Now, for each $j$, the quantizer $Q_j$ quantizes points in $R_j$, the range space of $P_j$. We can assume that for all $z \in R_j$, $Q_j(z) \in R_j$. Therefore, for all $d$, $P_j(d) - Q_j(P_j(d)) \in R_j$, and hence $P_j(d) - Q_j(P_j(d))$ is orthogonal to $d - P_T(d)$. Therefore,

$$\|d - Q(d)\|^2 = \|d - P_T(d)\|^2 + \|P_T(d) - Q_T(P_T(d))\|^2.$$

Using (17),

$$N D_{\min} = \mathbf{E}\|d - Q(d)\|^2 - \mathbf{E}\|P_T(d) - Q_T(P_T(d))\|^2. \qquad (19)$$

We now bound the two terms on the right hand side of (19).

For the first term, the quantized point $Q(d)$ can be described by $\log_2 J$ bits to quantize the index $T$ plus $n$ bits for the point $Q_T(P_T(d))$. Therefore, $Q(d)$ can be described by a total of $n + \log_2 J$ bits. The term $\mathbf{E}\|d - Q(d)\|^2$ is the average distortion of the quantizer $Q$ on the source $d$. Since $d$ is an $N$-dimensional jointly Gaussian vector with covariance $\sigma^2 I$, the distortion is bounded below by the distortion–rate function [8] to give

$$\mathbf{E}\|d - Q(d)\|^2 \geq N\sigma^2 2^{-2(n+\log_2 J)/N}. \qquad (20)$$

For the second term on the right hand side of (19), let

$$\sigma_j^2 = \mathbf{E}\left(\|P_j(d)\|^2 \mid T = j\right).$$

Now, in general, the distribution of $P_j(d)$ conditional on $T = j$ is *not* Gaussian. However, the distortion achievable for any distribution is always less than or equal

16

to the minimum distortion for a white Gaussian source with the same total variance. Consequently, for every $j$, the quantizer $Q_j$ can attain a distortion for a $K$-dimensional white Gaussian source with total variance $\sigma_j^2$. Therefore,

$$\mathbf{E}\left(\|P_j(d) - Q_j(P_j(d))\|^2 \mid T = j\right) \leq \sigma_j^2 2^{-2n/K}. \tag{21}$$

Also, observe that

$$\mathbf{E}\sigma_T^2 = \mathbf{E}\|P_T(d)\|^2 = \mathbf{E}\|d\|^2 - \mathbf{E}\|d - P_T(d)\|^2 = N\sigma^2 - D. \tag{22}$$

Substituting (22) in (21),

$$\mathbf{E}\left(\|P_T(d) - Q_T(P_T(d))\|^2\right) \leq (N\sigma^2 - D)2^{-2n/K}. \tag{23}$$

Then, substituting (20) and (23) in (19) we obtain

$$ND_{\min} \geq N\sigma^2 \left(\frac{2^{-2(n+\log_2 J)/N} - 2^{-2n/K}}{1 - 2^{-2n/K}}\right). \tag{24}$$

Since this bound must be true for all $n$, one can maximize with respect to $n$ to obtain the strongest bound. This maximization is messy; however, maximizing the numerator is easier and gives almost as strong a bound. The numerator is maximized when

$$n = \frac{NK \log\left(\frac{N}{K} J^{2/N}\right)}{(2\log 2)(N - K)},$$

and substituting this value of $n$ in (24) gives

$$
\begin{aligned}
ND_{\min} &\geq N\sigma^2 \cdot J^{-1/r} \cdot \frac{\left(\frac{K}{N}\right)^{K/(N-K)} - \left(\frac{K}{N}\right)^{N/(N-K)}}{1 - \left(\frac{K}{N}\right)^{N/(N-K)}} \\
&> N\sigma^2 \cdot J^{-1/r} \cdot \left(\left(\frac{K}{N}\right)^{K/(N-K)} - \left(\frac{K}{N}\right)^{N/(N-K)}\right).
\end{aligned}
$$

Dividing by $N$ completes the proof.

# B  Proof of Theorem 2

The proof is a slight modification of the proof of Theorem 1. Fix a signal $x \in \mathbb{R}^N$. The estimate $\hat{x}_{ML}$ is the projection of the noisy signal $y$ onto the subspace with index $T$. Thus, conditional on $T$,

$$\mathbf{E}((y - \hat{x}_{ML})'\hat{x}_{ML} \mid T, x) = 0.$$

17

Taking the expectation over $T$, we obtain the orthogonality relationship

$$\mathbf{E}((y - \hat{x}_{ML})'\hat{x}_{ML} \mid x) = 0.$$

Also, since we have assumed that the estimator $\hat{x}_{ML}$ is unbiased and the noise $d$ has zero mean,

$$\mathbf{E}((y - \hat{x}_{ML})'x \mid x) = \mathbf{E}((x + d - \hat{x}_{ML})'x \mid x) = x'x - x'x = 0.$$

Combining these orthogonality relationships,

$$\mathbf{E}\left[(y - \hat{x}_{ML})'(x - \hat{x}_{ML}) \mid x\right] = 0,$$

and therefore,

$$\begin{aligned}
\mathbf{E}\left[\|x - \hat{x}_{ML}\|^2 \mid x\right] &= \mathbf{E}\left[\|y - x\|^2 \mid x\right] - \mathbf{E}\left[\|y - \hat{x}_{ML}\|^2 \mid x\right] \\
&= \mathbf{E}\left[\|d\|^2\right] - \mathbf{E}\left[\|y - \hat{x}_{ML}\|^2 \mid x\right] \\
&= N\sigma^2 - \mathbf{E}\left[\|y - \hat{x}_{ML}\|^2 \mid x\right],
\end{aligned} \tag{25}$$

so we need to evaluate the expected difference $\|y - \hat{x}_{ML}\|^2$. If we define

$$\hat{d} = x - \hat{x}_{ML}$$

then

$$y - \hat{x}_{ML} = x + d - \hat{x}_{ML} = d - \hat{d}.$$

Now, the expected error $\mathbf{E}\left[\|d - \hat{d}\|^2 \mid x\right]$ can be lower bounded using analysis similar to Theorem 1. Also, the bound can be tightened by replacing the $\log_2 J$ bits with the entropy $H(T)$. Performing these computations, we obtain that for large $N$ and small $K/N$,

$$\mathbf{E}\left[\|y - \hat{x}_{ML}\|^2 \mid x\right] = \mathbf{E}\left[\|d - \hat{d}\|^2 \mid x\right] \geq N\sigma^2 2^{-H(T)/r}(1 - K/N). \tag{26}$$

Substituting (26) into (25) shows

$$e(x) = \frac{1}{N}\mathbf{E}\left[\|x - \hat{x}_{ML}\|^2 \mid x\right] = \sigma^2\left(1 - 2^{-H(T)/r}(1 - K/N)\right).$$

# C  Proof of Theorem 3

We begin with two lemmas. The first describes the distance between a unit vector and a single, random $K$-dimensional subspace.

**Lemma 2** *Suppose $z \in \mathbb{R}^N$ with $\|z\| = 1$. Consider a subspace $V$ of $\mathbb{R}^N$ spanned by $K$ random vectors uniformly distributed on the unit sphere, and let $\rho \in [0,1]$ be the distance squared from $z$ to $V$. Then, for small $\epsilon$,*

$$\Pr(\rho < \epsilon) = \frac{1}{r\beta(r, K/2)} \epsilon^r$$

*where $r$ is defined in (16) above.*

**Proof:** Since the distribution of the random vectors is spherically symmetric, we can consider the subspace $V$ fixed and take $z$ to be a random vector uniformly distributed on the unit sphere. One way to create such a random vector $z$ is to take $z = w/\|w\|$, where $w$ is unit variance, zero-mean Gaussian white noise. Let $w_1, \ldots, w_K$ be the components of $w$ on $V$, and $w_{K+1}, \ldots, w_N$ be the components in the orthogonal complement to $V$. If we define

$$X = \sum_{i=1}^{K} w_i^2 \quad \text{and} \quad Y = \sum_{i=K+1}^{N} w_i^2,$$

then the distance squared from $z$ to $V$ is

$$\rho = Y/(X + Y).$$

Now, define

$$U = \frac{Y/(N-K)}{X/K}.$$

Thus, $\rho < \epsilon$ if and only if

$$U = \frac{Y}{aX} = \frac{\rho}{a(1-\rho)} < \frac{\epsilon}{a(1-\epsilon)},$$

where $a = (N-K)/K$. Since $Y$ and $X$ are the sums of $N - K$ and $K$ independent Gaussians, $U$ has the F distribution with $N - K$ and $K$ degrees of freedom [35]. Therefore $U$ has the probability density function

$$f_U(u) = \frac{a^r u^{r-1}}{\beta(r, K/2)(1 + au)^{N/2}}.$$

For small $u$, the probability density function simplifies to

$$f_U(u) \approx \frac{a^r}{\beta(r, K/2)} u^{r-1}.$$

19

Thus, for small $\epsilon$,

$$\begin{aligned}
\Pr(\rho < \epsilon) &= \Pr\left(U < \frac{\epsilon}{a(1-\epsilon)}\right) \approx \Pr(U < \epsilon/a) \\
&= \int_0^{\epsilon/a} f_U(u)du \approx \frac{1}{r\beta(r, K/2)}\epsilon^r.
\end{aligned}$$

The next lemma describes the minimum distance between a unit vector and the closest of the $J$ subspaces spanned by $K$ vectors in the random frame $\Phi$.

**Lemma 3** *Let $z \in \mathbb{R}^N$ with $\|z\| = 1$, and consider the random frame $\Phi$ in the statement of the theorem. Let $\rho_{\min}$ be the minimum distance squared from $z$ to all the subspaces spanned by $K$ of the $M$ vectors in $\Phi$. If $M \gg K$ and $\epsilon$ is small,*

$$\Pr(\rho_{\min} > \epsilon) = \exp\left(-(C\epsilon)^r\right),$$

*where $J$ and $C$ are defined in (16).*

**Proof:** There are $J$ subspaces spanned by $K$ of the $M$ vectors in $\Phi$. Let $\rho_j$ be the distance squared from $z$ to the subspace with index $j$ for $j = 1, 2, \ldots, J$, so $\rho_{\min}$ is the minimum of these distances. Since $M \gg K$, we can assume that distances $\rho_j$ are independent. Now, if we set

$$C_0 = \frac{1}{r\beta(r, K/2)}$$

and apply the previous lemma, we obtain

$$\begin{aligned}
\Pr(\rho_{\min} > \epsilon) &= \prod_{j=1}^{J} \Pr(\rho_j > \epsilon) \\
&= [1 - C_0\epsilon^r]^J \approx \exp\left(-C_0 J\epsilon^r\right) = \exp\left(-(C\epsilon)^r\right),
\end{aligned}$$

where we have used the fact that $J$ is large.

**Proof of Theorem:** Let $V_0$ be the true subspace. That is, $V_0$ is a subspace containing the true signal $x$ and spanned by $K$ of the frame vectors in $\Phi$. Let $D_0$ be the distance squared from the noisy signal $y$ to $V_0$, and let $D_{\min}$ be the minimum distance squared from $y$ to the closest of the $J$ subspaces spanned by $K$ vectors in $F$. The estimator will select the correct subspace when $D_{\min} > D_0$. Therefore,

$$p_{\text{corr}} = \Pr(D_{\min} > D_0). \tag{27}$$

To compute the probability in (27), let $d_0$ be the $K$-dimensional component of $d$ in $V_0$, and let $d_1$ be the $N - K$ dimensional component orthogonal to $V_0$. Therefore,

20

$D_0$, the distance squared from $y$ to $V_0$, is $\|d_1\|^2$. Since the noise $d$ has variance $\sigma^2$ per dimension,

$$D_0 = \|d_1\|^2 = \sigma^2 U$$

where $U$ is a $\chi^2$ variable with $N - K = 2r$ degrees of freedom. If we let $U_{2r} = U/(2r)$ be the normalized $\chi^2$ variable,

$$D_0 = (N - K)\sigma^2 U_{2r}. \tag{28}$$

Also, since $N \gg K$, $\|d_1\| \gg \|d_0\|$. Therefore,

$$
\begin{aligned}
\|y\|^2 &= \|x + d\|^2 = \|x_+ d_0\|^2 + \|d_1\|^2 \\
&\approx \|x\|^2 + \|d_1\|^2 = N + D_0.
\end{aligned}
$$

Let $z = y/\|y\|$ and $\rho_{\min}$ be the minimum distance from $z$ to the closest subspace spanned by $K$ vectors in the random frame $\Phi$. Since $\|z\| = 1$, the distribution of $\rho_{\min}$ is described by Lemma 3. Also, since $y = \|y\|z$, the squared distance, $D_{\min}$ from $y$ to the closest subspace is given by

$$D_{\min} = \|y\|^2 \rho_{\min} = (1 + D_0)\rho_{\min}.$$

Therefore, applying Lemma 3 and equations (27), (28) and (29),

$$
\begin{aligned}
p_{\text{corr}} &= \Pr(D_{\min} > D_0) = \Pr\left(\rho_{\min} > \frac{D_0}{1 + D_0}\right) \\
&= \Pr\left(\rho_{\min} > \frac{(N - K)\sigma^2 U_{2r}}{N + (N - K)\sigma^2 U_{2r}}\right) \\
&= \Pr\left(\rho_{\min} > \frac{a\sigma^2 U_{2r}}{1 + a\sigma^2 U_{2r}}\right) \\
&= E\exp\left(-\left(\frac{a\sigma^2 C U_{2r}}{1 + a\sigma^2 U_{2r}}\right)^r\right)
\end{aligned}
$$

Now, it can be verified using the formulae in [34] that the normalized $\chi^2$ variable $U_{2r}$ has a probability density function $p_r(u)$ in (15). It follows that

$$p_{\text{corr}} = \int_0^\infty p_r(u) \exp\left(-\left(\frac{a\sigma^2 C u}{1 + a\sigma^2 u}\right)^r\right) = \phi(a\sigma^2, C, r).$$

This completes the proof.

# References

[1] F. Bergeaud and S. Mallat. Matching pursuit of images. In *Proc. IEEE Int. Conf. Image Proc.*, volume I, pages 53–56, Washington, DC, October 1995.

[2] S. G. Chang, B. Yu, and M. Vetterli. Adaptive wavelet thresholding for image denoising and compression. *IEEE Trans. Image Proc.*, 9(9):1532–1546, September 2000.

[3] S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM J. Sci. Comp.*, 20(1):33–61, 1999.

[4] S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM Rev.*, 43(1):129–159, 2001.

[5] A. Cohen and J.-P. D'Ales. Nonlinear approximation of random functions. *SIAM J. Appl. Math.*, 57(2):518–540, April 1997.

[6] J. H. Conway, R. H. Hardin, and N. J. A. Sloane. Packing lines, planes, etc.: Packings in Grassmannian spaces. *Experimental Mathematics*, 5(2):139–159, 1996. See also [7].

[7] Editors' note on packing lines, planes, etc.: Packings in Grassmannian spaces. *Experimental Mathematics*, 6(2):175, 1997.

[8] T. M. Cover and J. A. Thomas. *Elements of Information Theory*. John Wiley & Sons, New York, 1991.

[9] G. Davis. *Adaptive Nonlinear Approximations*. PhD thesis, New York Univ., September 1994.

[10] R. A. DeVore. Nonlinear approximation. *Acta Numerica*, pages 51–150, 1998.

[11] D. Donoho and M. Elad. Optimally sparse representation in general (nonorthogonal) dictionaries via $\ell^1$ minimization. *Proc. Nat. Acad. Sci.*, 100(5):2197–2202, March 2003.

[12] D. L. Donoho, M. Elad, and V. Temlyakov. Stable recovery of sparse overcomplete representations in the presence of noise. 2004.

[13] D. L. Donoho, M. Vetterli, R. A. DeVore, and I. Daubechies. Data compression and harmonic analysis. *IEEE Trans. Inform. Th.*, 44(6):2435–2476, October 1998.

[14] M. Elad and A. M. Bruckstein. A generalized uncertainty principle and sparse representation in pairs of bases. *IEEE Trans. Inform. Th.*, 48(9):2558–2567, September 2002.

[15] K. Engan, S. O. Aase, and J. H. Husøy. Designing frames for matching pursuit algorithms. In *Proc. IEEE Int. Conf. Acoust., Speech, and Signal Proc.*, volume 3, pages 1817–1820, Seattle, WA, May 1998.

[16] A. K. Fletcher and K. Ramchandran. Estimation error bounds for frame denoising. In *Proc. Wavelets: Appl. in Sig. & Image Proc. X, part of SPIE Int. Symp. on Optical Sci. & Tech.*, volume 5207, pages 40–46, San Diego, CA, August 2003.

[17] G. H. Golub and C. F. Van Loan. *Matrix Computations*. Johns Hopkins Univ. Press, Baltimore, MD, second edition, 1989.

[18] M. M. Goodwin. *Adaptive Signal Models: Theory, Algorithms and Audio Applications*. Kluwer Acad. Pub., 1998.

[19] V. K Goyal, M. Vetterli, and N. T. Thao. Quantized overcomplete expansions in $\mathbb{R}^N$: Analysis, synthesis, and algorithms. *IEEE Trans. Inform. Th.*, 44(1):16–31, January 1998.

[20] R. Gribonval and M. Nielsen. Sparse representations in unions of bases. *IEEE Trans. Inform. Th.*, 49(12):3320–3325, December 2003.

[21] H. Krim, D. Tucker, S. Mallat, and D. Donoho. On denoising and best signal representation. *IEEE Trans. Inform. Th.*, 45(7):2225–2238, November 1999.

[22] J. Liu and P. Moulin. Complexity-regularized image denoising. *IEEE Trans. Image Proc.*, 10(6):841–851, June 2001.

[23] S. G. Mallat and Z. Zhang. Matching pursuits with time-frequency dictionaries. *IEEE Trans. Signal Proc.*, 41(12):3397–3415, December 1993.

[24] F. Moschetti, L. Granai, P. Vandergheynst, and P. Frossard. New dictionary and fast atom searching method for matching pursuit representation of displaced frame difference. In *Proc. IEEE Int. Conf. Image Proc.*, volume 3, pages 685–688, Rochester, NY, September 2002.

[25] B. K. Natarajan. Filtering random noise from deterministic signals via data compression. *IEEE Trans. Signal Proc.*, 43(11):2595–2605, November 1995.

[26] B. K. Natarajan. Sparse approximate solutions to linear systems. *SIAM J. Computing*, 24(2):227–234, April 1995.

[27] R. Neff and A. Zakhor. Very low bit-rate video coding based on matching pursuits. *IEEE Trans. Circuits Syst. Video Technol.*, 7(1):158–171, February 1997.

[28] N. Saito. Simultaneous noise suppression and signal compression using a library of orthonormal bases and the minimum description length criterion. In E. Foufoula-Georgiou and P. Kumar, editors, *Wavelets in Geophysics*, pages 299–324. Academic Press, San Diego, CA, 1994.

[29] N. J. A. Sloane, R. H. Hardin, and W. D. Smith. A library of putatively optimal spherical codes, together with other arrangements which may not be optimal but are especially interesting for some reason. URL: `http://www.research.att.com/~njas/packings.` ·

[30] T. Strohmer and R. W. Heath Jr. Grassmannian frames with applications to coding and communication. *Appl. Comput. Harm. Anal.*, 14(3):257–275, May 2003.

[31] J. A. Tropp. Greed is good: Algorithmic results for sparse approximation. ICES Report 0304, Univ. of Texas at Austin, February 2003.

[32] J. A. Tropp. Just relax: Convex programming methods for subset selection and sparse approximation. ICES Report 0404, Univ. of Texas at Austin, February 2004.

[33] E. W. Weisstein. Beta function. From *MathWorld*–A Wolfram Web Resource, http://mathworld.wolfram.com/BetaFunction.html.

[34] E. W. Weisstein. Chi-squared distribution. From *MathWorld*–A Wolfram Web Resource, http://mathworld.wolfram.com/Chi-SquaredDistribution.html.

[35] E. W. Weisstein. *F*-distribution. From *MathWorld*–A Wolfram Web Resource, http://mathworld.wolfram.com/F-Distribution.html.

[36] E. W. Weisstein. Gamma function. From *MathWorld*–A Wolfram Web Resource, http://mathworld.wolfram.com/GammaFunction.html.