

## Chapter 2

# How (Well) Compressed Sensing Works in Practice

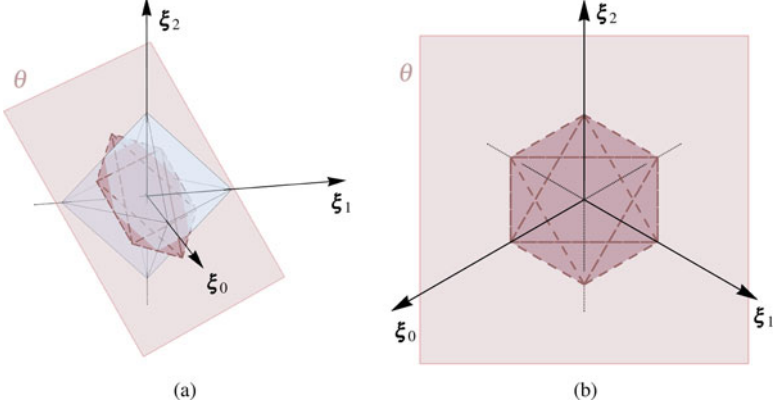
### 2.1 Non-Worst-Case Assessment of CS Performance

One of the main problems with coherence and restricted isometries is that the corresponding parameters are explicitly calibrated on worst-case scenarios. This corresponds to the desire of providing guarantees, thanks to which the system is known to operate correctly. Yet, acquisition systems work on random inputs and it is perfectly sensible to characterize their performance by probabilistic means. This is particularly true when the input is not their unique random components since, for example, the matrix  $A$  is also a possibly time varying, uncertain ingredient of the processing.

An instructive example of this alternative route is given by a more geometric approach to the properties of the minimization problem (1.27) (Basis Pursuit—BP) and to its relationship with the minimization problem (1.26). Everything hinges on a special kind of polytopes.

**Definition 2.1** These are the definitions we need to proceed in our analysis

- A  $p$ -dimensional convex polytope  $P \subset \mathbb{R}^p$  is the convex hull of a set  $V$  of points in  $\mathbb{R}^p$ .
- If no point in  $V$  can be dropped without changing the resulting convex hull, then the points in  $V$  are the *vertices* of  $P$ .
- The intersection  $P \cap h$  of a  $p$ -dimensional convex polytope  $P$  with a  $p - 1$ -dimensional hyperplane that does not contain any point of the interior of  $P$  is called a *facet* of  $P$ . Facets can be 0-dimensional (vertices), 1-dimensional (edges), or in general  $q$ -dimensional with  $q < p$ .
- Given a convex polytope  $P \subset \mathbb{R}^p$  with vertices  $\mathbf{v}_0, \mathbf{v}_1, \dots$  and a  $q \times p$  matrix  $\mathbf{B}$ , the convex hull of the points  $\mathbf{B}\mathbf{v}_0, \mathbf{B}\mathbf{v}_1, \dots$  is a  $q$ -dimensional convex polytope  $Q = \mathbf{B}P \subset \mathbb{R}^q$ .



**Fig. 2.1** Construction of  $BS_3^1(1)$  starting from  $S_3^1(1)$  using  $B$  as in (1.7) (a) and its frontal view allowing face counting (b)

- A polytope is said to be *centrosymmetric* if  $\mathbf{v} \in V$  implies  $-\mathbf{v} \in V$ .
- If  $\mathbf{e}_j = (0, \dots, 0, 1, 0, \dots, 0)^T$  where the unique 1 appears in the  $j$ -th position, then the *crosspolytope* in  $\mathbb{R}^p$  is defined as the centrosymmetric convex hull of  $V = \{\pm \mathbf{e}_0, \pm \mathbf{e}_1, \dots, \pm \mathbf{e}_{p-1}\}$ . In our previous notation the crosspolytope is nothing but  $S_p^1(1)$ .

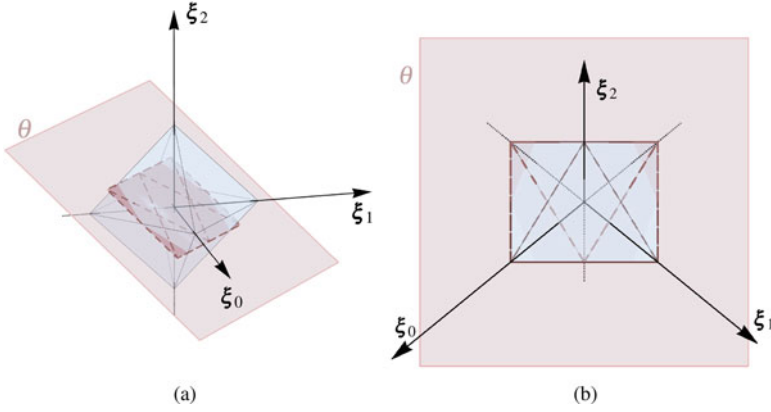
Starting from these definitions that are standard elements in the theory of convex polytopes, one may develop specific concepts related to possibility of reconstructing the original signal from the measurements vector [7, 8, Theorem 7.5].

**Definition 2.2** A  $p$ -dimensional centrosymmetric polytope is said to be *centrally  $q$ -neighborly* if every subset of  $V$  with  $q$  elements that does not include two antipodal vertices is the set of vertices of a  $q - 1$ -dimensional facet of  $P$ .

**Theorem 2.1** If  $\mathbf{y} = B\boldsymbol{\xi}$  has a unique solution with not more than  $\kappa$  non-null components, then such a solution is the unique solution of  $BP$  if and only if  $BS_d^1(1)$  has  $2d$  vertices and is centrally  $(\kappa - 1)$ -neighborly.

The point of interest in Theorem 2.1 is that it gives a necessary and sufficient condition for signal reconstruction: no worst-case bounding is involved. This has a substantial impact on the predictability of CS performance. As an example, we may go back to the matrix  $B$  in (1.7) and recall that its mutual coherence equal to 1 and its RIC equal to  $1/2$  prevent the application of the results leveraging those concepts, i.e., of Theorems 1.5, 1.6, and 1.7.

Yet, Theorem 2.1 explains why reconstruction by means of BP is always effective in the noiseless case. Figure 2.1a reports the construction of  $BS_3^1(1)$  starting from  $S_3^1(1)$ . The same  $BS_3^1(1)$  is visualized in Fig. 2.1b. In this case  $d = 3$  and  $\kappa = 1$  and it is easy to verify that  $BS_3^1(1)$  has  $2d = 6$  vertices each of them trivially being a face so that the polytope is centrally 0-neighborly.



**Fig. 2.2** Construction of  $BS_3^1(1)$  starting from  $S_3^1(1)$  using  $B$  as in (2.1) (a) and its frontal view allowing face counting (b)

The same theorem indicates when the choice of  $B$  may prevent us from getting a reconstruction. As an example consider

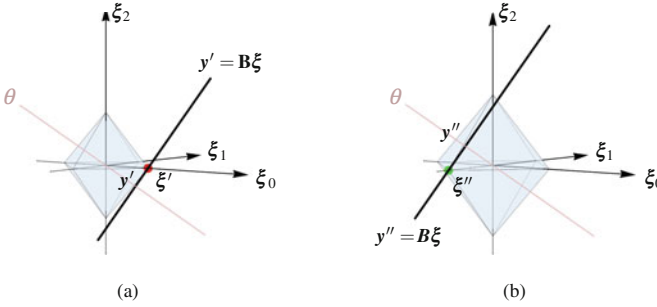
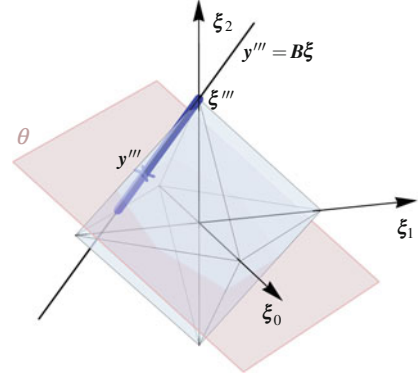
$$B = \begin{pmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & 0 \\ -\frac{1}{\sqrt{3}} & -\frac{1}{\sqrt{3}} & \frac{1}{\sqrt{3}} \end{pmatrix} \quad (2.1)$$

and the resulting  $BS_3^1(1)$  as in Fig. 2.2. In this case, the number of vertices is only  $4 < 2d = 6$  and Theorem 2.1 implies that reconstruction by means of (1.27) may be impossible since the sparsity prior may not be enough to select a unique solution of  $y = B\xi$ . What happens is described in Fig. 2.2 in which  $S_3^1(1)$  is projected with the new  $B$  on a new plane  $\theta$ . Two out of 8 faces of  $S_3^1(1)$  are orthogonal to  $\theta$  so that one of the vertices of each of these two faces gets mapped into a point on the projection of the edge connecting the other two vertices, disappearing in the final polytope  $BS_3^1(1)$ .

This is what prevents reconstruction. In fact, assume that you want to recover the same point  $\xi'''$  as in Fig. 1.14a. The straight line corresponding to  $y''' = B\xi$  that is orthogonal to  $\theta$  is also parallel to 2 faces of  $S_3^1(1)$  so that its intersection with  $S_3^1(\|\xi'''\|_1)$  is a whole segment, each point of which is a solution of BP. This is visualized in Fig. 2.3.

As a further application to our toy case, the very same Theorem 2.1 explains why, no matter how the plane  $\theta$  is positioned, the solution of BP is not able to retrieve  $\xi$  from  $y = B\xi$  when it is known that  $\xi$  is 2-sparse instead of 1-sparse. In fact, to have  $2d$  vertices,  $BS_3^1(1)$  must be a hexagon. Yet, only the pairs of consecutive vertices belong to a  $(\kappa - 1)$ -dimensional facet (that for  $\kappa = 2$  is an edge) of a hexagon, while other pairs do not, preventing central neighborliness. Hence, no matter how  $\theta$  is positioned, BP cannot be used for signal reconstruction.

**Fig. 2.3** Trying the reconstruction of  $\xi'''$  starting from  $y''' = B\xi'''$  by means of BP. All the points in the blue thick segment are equally good solutions



**Fig. 2.4** Successful reconstruction of  $\xi'$  starting from  $y' = B\xi'$  (a) and of  $\xi''$  from  $y'' = B\xi''$  (b) by means of BP

Lastly, this powerful point of view can be extended to the case in which the signal to retrieve is random. In fact, Theorem 2.1 is a guarantee independent of  $\xi$ , that ceases to hold if even a single  $\xi$  cannot be reconstructed. Yet, even when the guarantee does not hold, like for the matrix in (2.1), there are signals that can be reconstructed.

In particular, in passing from  $S_3^1(1)$  to  $BS_3^1(1)$ , 2 out of 6 vertices are lost and the signals  $\xi$  that cannot be retrieved are exactly those at the corresponding vertices of  $S_3^1(\|\xi\|_1)$ , like  $\xi'''$  in Fig. 2.3. Yet, Fig. 2.2 shows that two other pairs of vertices appear in  $BS_3^1(1)$  and signals on the corresponding vertices of  $S_3^1(\|\xi\|_1)$  can still be reconstructed. This is shown in Figs. 2.4a and b where the same  $\xi'$  and  $\xi''$  as in Fig. 1.6 are uniquely identified by the intersection of the straight line  $y = B\xi$  and the minimum radius  $\|\cdot\|_1$  ball.

Intuitively speaking, if the original 1-sparse signal  $\xi$  has the same probability of aligning with each of the axes, the probability that BP is effective in recovering it is equal to the ratio of the number of surviving vertices over the number of original vertices, i.e.,  $4/6 = 2/3$ .

All this can be generalized to cope with a larger sparsity  $\kappa$ . To understand how, we may first define  $\phi_\kappa(\cdot)$  as the operator that counts the number of  $\kappa$ -dimensional facets of its polytope argument and state the following [7, Theorem 3].

**Theorem 2.2** *Let  $\mathbf{B}$  be an  $m \times d$  matrix such that if  $\mathbf{B}\boldsymbol{\alpha} = 0$  for a vector  $\boldsymbol{\alpha}$  with less than  $m$  nonzeros then  $\boldsymbol{\alpha} = 0$ . Let also  $\kappa < m/2$ .*

*Given a subset  $K \subset \{0, \dots, d-1\}$  of cardinality  $\kappa$ , we may have that if  $\text{supp}(\boldsymbol{\xi}) = K$  then  $\boldsymbol{\xi}$  can be reconstructed from  $\mathbf{y} = \mathbf{B}\boldsymbol{\xi}$  by means of BP. Indicate with  $K_{\text{BP}}$  the number of such subsets, and with  $K_{\text{tot}} = \binom{d}{\kappa}$  the total number of possible subsets of cardinality  $\kappa$ . Then*

$$\frac{K_{\text{BP}}}{K_{\text{tot}}} \geq \frac{\phi_{\kappa-1}(\mathbf{B}S_d^1(1))}{\phi_{\kappa-1}(S_d^1(1))} \quad (2.2)$$

Assuming that the original signal has the same probability of featuring any of the  $K_{\text{tot}}$  supports of cardinality  $\kappa$ , the above result can be immediately recast into probabilistic terms to say that  $p_{\text{BP}} = K_{\text{BP}}/K_{\text{tot}}$  is the probability of successful reconstruction by means of BP and is not less than the ratio of facets counts in (2.2). A dual result is available for the case in which  $\mathbf{B}$  is random [8, Theorem 7.7].

**Theorem 2.3** *Let  $\mathbf{B}$  be a random  $m \times d$  matrix whose probability distribution is invariant for any signed permutation of rows. Let  $\boldsymbol{\xi} \in \mathbb{R}^d$  be a  $\kappa$ -sparse vector and  $\mathbf{y} = \mathbf{B}\boldsymbol{\xi}$  the corresponding random measurement vector. The probability  $p_{\text{BP}}$  that BP retrieves  $\boldsymbol{\xi}$  from  $\mathbf{y}$  is bounded by*

$$p_{\text{BP}} \geq \frac{\mathbb{E}[\phi_{\kappa-1}(\mathbf{B}S_d^1(1))]}{\phi_{\kappa-1}(S_d^1(1))}$$

Note that, in general, facets counting is a combinatorial task so that the computation of  $\phi_{\kappa-1}(\cdot)$  in high-dimensional settings can be expensive if not impossible. From this point of view, the introduction of random matrices  $\mathbf{B}$  can be helpful if paired with the asymptotic conditions that are the mathematical equivalent of the high-dimensional setting in which CS is applied. Many sophisticated results are born in this area, whose simplest prototype is probably the one that we rephrase here in our terms [10].

**Theorem 2.4** *Let  $\mathbf{B} \sim \text{RGE}(\text{iid})$  with unit variance entries,  $d = (\text{DR} \times \text{CR})m$ , and  $m = \text{OH } \kappa$ . There is a function  $\psi(\cdot)$  such that*

$$\lim_{d \rightarrow \infty} \frac{\phi_{\kappa-1}(\mathbf{B}S_d^1(1))}{\phi_{\kappa-1}(S_d^1(1))} = \begin{cases} 1 & \text{if } \text{DR} \times \text{CR} < \psi(\text{OH}) \\ 0 & \text{if } \text{DR} \times \text{CR} > \psi(\text{OH}) \end{cases}$$

Collecting the results in Theorems 2.2, 2.3, and 2.4 one gets that, as the dimensionality increases, there is a crisp *phase transition* in the possibility of

reconstructing  $\xi$  from its random projections. The excess of measurements with respect to the actual degrees of freedom in the signal (OH) controls the possibility of accommodating a certain dimensionality reduction  $DR \times CR$  while maintaining the retrievability of the original signal.

Though Theorem 2.4 leverages RGE (iid), the existence and shape of the function  $\psi$  has been empirically found to be a general property [9] when the entries of  $\mathbf{B}$  are iid or its rows are an iid random subset of certain orthonormal basis.

In the noiseless and exactly sparse case, this makes polytope-based analysis much closer to real performance than coherence or RIP-based considerations since neither the finite-dimension results, nor asymptotics of random  $\mathbf{B}$  rely on worst-case bounding.

As a consequence, though neither the theory which heavily relies on symmetry considerations, nor the empirical evidence gathered so far in the Literature, say much on the possibility of straightforwardly applying this point of view to the design of a proper sensing matrix  $\mathbf{A}$ , we know that if  $n$  is large and we increase  $m$  enough, CS will eventually work very well (the *probability 1* implicit in Theorem 2.4).

From a more engineering point of view, this can be reversed to say that our aim is to find the minimum possible  $m$  for which CS works very well. Properly evolved and specialized, this is the key idea behind the discussion in the chapters to follow.

## 2.2 Beyond Basis Pursuit

Despite its theoretical appeal, BP is only an archetypal reconstruction method. In practical terms BP and its denoising variant BPDN have been implemented with a variety of methods, ranging from straightforward mapping to classical mathematical programming problems leveraging linear and quadratic optimization tools, to specialized procedures that look at them as a particular case of a convex optimization task.

The activity in this field revealed, for example, that it is sometimes convenient to address the BP or BPDN problems not in the *synthesis form* contained in (1.27) but in an alternative *analysis form*.

Note, in fact, that (1.27) depends only on  $\mathbf{B}$  and thus considers and seeks to reconstruct the signal  $\xi$  in the sparsity domain. Once that  $\xi$  is known one may *synthesize*  $\mathbf{x} = \mathbf{D}\xi$ .

Assume now that a linear operator  $\mathbf{D}^*$  is available such that if  $\xi = \mathbf{D}^*\mathbf{x}$  then  $\mathbf{x} = \mathbf{D}\xi$ . When  $\mathbf{D}$  is a non-singular square matrix we simply have  $\mathbf{D}^* = \mathbf{D}^{-1}$  and when  $\mathbf{D}$  is a frame,  $\mathbf{D}^*$  is the dual frame operator. With this, we may concentrate directly on the true signal  $\mathbf{x}$  and try to solve the “equivalent”

$$\begin{aligned} \arg \min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{D}^*\mathbf{x}\|_1 \\ \text{s.t. } \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2 \leq \epsilon \end{aligned} \tag{2.3}$$

Clearly, (1.27) and (2.3) are not equivalent when  $\mathbf{D}$  is not an invertible matrix. In fact,  $\mathbf{D}^*\boldsymbol{\xi}$  is only one of the many possible representations of  $\mathbf{x}$  in the sparsity domain and is the only one considered while scanning the feasibility space of (2.3) while (1.27) considers all of them. What happens is that the choice made by  $\mathbf{D}^*$  acts as a further prior and in this role, it is often useful to decreased dimensionality of the analysis form of BP and BPDN and help them finding good solutions.

Beyond this, accounting for all the methods and implementations described in the Literature and/or made available to practitioners is out of the scope of this book. Yet, it is useful to mention some of the most widespread tools distinguishing between those helping the implementation of BP, BPDN, and their other variants, those tackling the reconstruction problem from a theoretically different point of view, and those that are mainly based on heuristic considerations and yield lightweight iterative procedures that may be extremely useful when the resources dedicated to signal retrieval are limited.

Further to their implementation in commercial, large-scale solvers, BP and BPDN can be solved by quite a few implementations. Among them it is worthwhile mentioning those in Table 2.1 where we give the commonly used acronym, a pointer to some ready-to-use code and references to the relevant Literature.

Since BP and BPDN are convex optimization problems, they can be tackled by convex solvers with wider applicability. Those in Table 2.2 are particularly effective in modeling and solving the two standard reconstruction problems. Additionally, their greater generality can be used to add constraints that model priors further to sparsity that may available on the signal, thus increasing reconstruction performance.

Further to these methods, instead of depending on the  $\|\cdot\|_1$  norm and its favorable geometry, signal reconstruction can be approached from completely different points of view, e.g., from the estimation, or machine learning, or regression point of view. Different approaches result in different algorithms some of which are listed in Table 2.3.

**Table 2.1** Some dedicated BP/BPDN solvers

Solver	Url	Reference
SPGL1	<a href="http://www.math.ucdavis.edu/~mpf/spgl1/">www.math.ucdavis.edu/~mpf/spgl1/</a>	[2]
NESTA	<a href="http://statweb.stanford.edu/~candes/nesta/">statweb.stanford.edu/~candes/nesta/</a>	[1]

**Table 2.2** Some solvers of convex optimization problems that can be used for signal retrieval

Solver	Url	Reference
CVX	<a href="http://cvxr.com/">cvxr.com/</a>	[12, 13]
Unlocbox	<a href="http://lts2.epfl.ch/unlocbox/">lts2.epfl.ch/unlocbox/</a>	[5]

**Table 2.3** Some signal reconstruction methods based on various heuristic

Solver	Url	Reference
GAMP	<a href="http://gampmatlab.wikia.com">gampmatlab.wikia.com</a>	[17]
IRLS	<a href="http://stemblab.github.io/irls/">http://stemblab.github.io/irls/</a>	[6]
SBL	<a href="http://dsp.ucsd.edu/~zhilin/BSBL.html">dsp.ucsd.edu/~zhilin/BSBL.html</a>	[14]

**Table 2.4** Some signal reconstruction methods based on various heuristic

Solver	Url	Reference
FOCUSS	<a href="http://dsp.ucsd.edu/~jfmurray/software.htm">dsp.ucsd.edu/~jfmurray/software.htm</a>	[11]
OMP	<a href="http://www.mathworks.com/matlabcentral/fileexchange/32402-cosamp-and-omp-for-sparse-recovery">http://www.mathworks.com/matlabcentral/fileexchange/32402-cosamp-and-omp-for-sparse-recovery</a>	[16]
CoSaMP		[16]
Iterative hard thresholding	<a href="http://www.personal.soton.ac.uk/tb1m08/sparsify/sparsify.html">www.personal.soton.ac.uk/tb1m08/sparsify/sparsify.html</a>	[3]
	<a href="http://sparselab.stanford.edu/">http://sparselab.stanford.edu/</a>	

**Table 2.5** Code sketch for CoSaMP

<b>Require:</b> $\mathbf{y}$ vector of measurements	
<b>Require:</b> $\kappa$ sparsity level	
<b>Require:</b> $\mathbf{B} = \mathbf{A}\mathbf{D}$ sensing matrix	
$\hat{\xi} \leftarrow 0$	▷ signal guess
$\Delta\mathbf{y} \leftarrow \mathbf{y} - \mathbf{B}\hat{\xi} = \mathbf{y}$	▷ error in reproducing measurements
<b>repeat</b>	
$\Delta\hat{\xi} \leftarrow \mathbf{B}^\top \Delta\mathbf{y}$	▷ error in signal guess
$J = \text{supp}(\hat{\xi}) \cup \text{supp}(\Delta\hat{\xi}^{2\kappa\uparrow})$	▷ support to correct error in signal guess
$\hat{\xi} \leftarrow 0$	
$\hat{\xi}_J = (\mathbf{B}_{\cdot,J})^\dagger \mathbf{y}$	
$\hat{\xi} \leftarrow \hat{\xi}^{\kappa\uparrow}$	▷ new signal guess
$\Delta\mathbf{y} = \mathbf{y} - \mathbf{B}\hat{\xi}$	▷ new error in reproducing measurements
<b>until</b> convergence	

Finally, procedures exist that retrieve the original signal by considering that the main issue in the computation of  $\xi$  is not finding a generic solution to  $\mathbf{y} = \mathbf{B}\xi$  but to find the sparse one. Starting from this, it is possible to generate solutions iteratively adjusting their sparsity at each step. Different heuristics may be used to promote sparsity and give raise to different methods, some of which are listed in Table 2.4. The simple structure of these methods and their relatively good performance make them ideal for CS embodiments in which the resources devoted to signal reconstruction are limited.

As an example of how simple such algorithms can be, assume that  $\mathbf{B}$  is well approximated by a random matrix with i.i.d., zero-average, entries and that the  $\|\cdot\|_2$  norm of each column is approximately equal. Since the columns of  $\mathbf{B}$  are independent, the matrix  $\mathbf{B}^\top \mathbf{B}$  is well approximated by a diagonal matrix.

Assume now that an estimate  $\hat{\xi}$  is given of the true  $\xi$ . The measurement vector corresponding to  $\hat{\xi}$  is  $\hat{\mathbf{y}} = \mathbf{B}\hat{\xi}$  whose difference with respect to the true measurement vector is  $\Delta\mathbf{y} = \mathbf{B}(\hat{\xi} - \xi)$ . Thanks to the previous considerations on  $\mathbf{B}^\top \mathbf{B}$ , we also have that  $\mathbf{B}^\top \Delta\mathbf{y} = \mathbf{B}^\top \mathbf{B}(\hat{\xi} - \xi) \approx \|\mathbf{B}_{\cdot,0}\|_2 (\hat{\xi} - \xi)$ .

Hence, the largest nonzero components of  $\mathbf{B}^\top \Delta\mathbf{y}$  indicate the components of the signal that have been mistaken most by taking  $\hat{\xi}$  instead of  $\xi$ . This is the core step in the CoSaMP algorithm whose complete definition is given in Table 2.5 where:  $\cdot^{\uparrow p}$  that takes a vector and gives its thresholded version in which all but the  $p$



largest component are set to zero,  $\cdot^\dagger$  indicates the Moore–Penrose pseudo-inverse of a matrix, and given an index set  $J$ , a vector  $\mathbf{v}$ , and a matrix  $\mathbf{M}$ ,  $\mathbf{v}_J$  is the subvector of  $\mathbf{v}$  containing only the entries of  $\mathbf{v}$  with indexes in  $J$ , while  $\mathbf{M}_{\cdot,J}$  is the submatrix of  $\mathbf{M}$  made of the columns of  $\mathbf{M}$  whose indexes stay in  $J$ .

Though the convergence criterion is not specified, it is clear that the procedure itself is much simpler than solving a convex optimization problem. This is the reason why methods like this and like the others in Table 2.4 are often used in limited-resources realizations of reconstruction stages (see, e.g., [4]).

## 2.3 A Framework for Performance Evaluation

In the light of the discussions in Chap. 1 and of the initial section of this chapter, it is easy to state that a precise assessment of the performance of a CS system is far from easy.

An obvious intuition is that performance must be related to the magnitude of the reconstruction error, i.e., to the difference between the true sparse representation  $\xi$  and the one estimated by the reconstruction algorithm  $\hat{\xi}$  or between the true signal  $\mathbf{x}$  and  $\hat{\mathbf{x}} = \mathbf{D}\hat{\xi}$ .

Yet, the classical theory of Chap. 1 follows a worst-case leitmotif and gives bounds on quantities like  $\|\hat{\xi} - \xi\|_2$  that are either rarely applicable (for example, because they pose too strict requirements on measurement matrices  $\mathbf{A}$ ) or quite loose and ultimately very far from actual behavior.

Even the non-worst-case approach described at the beginning of this chapter has problems since its face-counting argument, though allowing a much sharper distinction between what can be reconstruct and what cannot, scales poorly as dimension increase and cannot be applied in practice.

Last but not least, the construction of the matrices  $\mathbf{A}$  is often done by random means. This, paired with the intrinsic random nature of the signal to acquire, implies that reconstruction error is a quite complicated random quantity.

The most straightforward way of addressing all these problems is to resort to extensive Montecarlo simulations. Such an approach is the most common both in the Literature and in practice and consists in generating a large number  $W$  of signal instances  $\mathbf{x}^{(j)}$  and of measurement matrices  $\mathbf{A}^{(j)}$  for  $j = 0, \dots, W - 1$ , use each of them to compute  $\mathbf{y}^{(j)} = \mathbf{A}^{(j)}\mathbf{x}^{(j)}$  and then run one of the algorithms mentioned before to compute the estimation  $\hat{\mathbf{x}}^{(j)}$  and consequently the reconstruction error in that case. The statistic of such an error is usually summarized in single numbers by means of one of the two approaches.

To begin with, it is most natural to define a Reconstruction Signal-to-Noise-Ratio

$$\text{RSNR}[\text{dB}] = 20 \log_{10} \left( \frac{\|\mathbf{x}\|_2}{\|\hat{\mathbf{x}} - \mathbf{x}\|_2} \right)$$

that acts as a merit figure, i.e., the larger the  $\text{RSNR}[\text{dB}]$ , the better the reconstruction.

Then one may try to estimate the Average RSNR[dB] as

$$\text{ARSNR[dB]} = \mathbf{E} \left[ 20 \log_{10} \left( \frac{\|\mathbf{x}\|_2}{\|\hat{\mathbf{x}} - \mathbf{x}\|_2} \right) \right] \approx \frac{1}{W} \sum_{j=0}^{W-1} 20 \log_{10} \left( \frac{\|\mathbf{x}^{(j)}\|_2}{\|\hat{\mathbf{x}}^{(j)} - \mathbf{x}^{(j)}\|_2} \right) \quad (2.4)$$

Alternatively, one may assume that the reconstruction is correct when the corresponding RSNR[dB] exceeds a certain  $\text{RSNR[dB]}_{\min}$  and define a Probability of Correct Reconstruction (PCR) as

$$\text{PCR} = \Pr\{\text{RSNR[dB]} \geq \text{RSNR[dB]}_{\min}\} \approx \frac{\left| \left\{ \frac{\|\mathbf{x}^{(j)}\|_2}{\|\hat{\mathbf{x}}^{(j)} - \mathbf{x}^{(j)}\|_2} \geq 10^{\frac{\text{RSNR[dB]}_{\min}}{20}} \right\} \right|}{W}$$

Clearly, ARSNR[dB] and PCR are general-purpose merit figures and real-world applications may provide more significant indexes for establishing the acquisition performance. When applications are addressed at the end of this book, those merit figures will be possibly described and applied.

Yet, examples made to describe the adaptive method we address will use ARSNR[dB] and PCR and a uniform framework for the accumulation of Montecarlo trials.

In particular, we are interested in  $n$ -dimensional signals  $\mathbf{x}$  that are both localized and  $\kappa$ -sparse with respect to a certain reference system  $\mathbf{D}$  that we assume to be an orthonormal basis.

To generate samples of  $\mathbf{x}$  we start from an instance of a zero-mean Gaussian random vector  $\mathbf{x}'$  with covariance/correlation matrix  $\mathcal{X}'$  and do the following steps:

$$\begin{aligned} \mathbf{x}' &\sim \mathbf{N}(0, \mathcal{X}') \\ \xi' &\leftarrow \mathbf{D}^{-1} \mathbf{x}' = \mathbf{D}^{\top} \mathbf{x}' \\ \xi &\leftarrow (\xi')^{\kappa \uparrow} \\ \mathbf{x} &= \mathbf{D} \xi \end{aligned}$$

that formalize the intuitive idea of taking a possibly non-white vector ( $\mathbf{x}'$ ) project it onto the basis along which we want our signal to be sparse, sparsify it and map it back into its original basis. Clearly, if  $\kappa = n$  we have  $\mathbf{x} = \mathbf{x}'$  since no clipping takes place.

As a first remark, note that from  $\mathbf{E}[\mathbf{x}'] = 0$  we have  $\mathbf{E}[\xi'] = 0$ ,  $\mathbf{E}[\xi] = 0$  and  $\mathbf{E}[\mathbf{x}] = 0$ . Moreover, if we define the covariance/correlation matrices  $\mathcal{X}' = \mathbf{E}[\mathbf{x}\mathbf{x}'^{\top}]$ ,  $\mathcal{E}' = \mathbf{E}[\xi'\xi'^{\top}]$ ,  $\mathcal{E} = \mathbf{E}[\xi\xi^{\top}]$ ,  $\mathcal{X} = \mathbf{E}[\mathbf{x}\mathbf{x}^{\top}]$ , we have  $\mathcal{E}' = \mathbf{D}^{\top} \mathcal{X}' \mathbf{D}$  and  $\mathcal{X} = \mathbf{D} \mathcal{E} \mathbf{D}^{\top}$ .

Hence, if  $\mathcal{X}'$  is a diagonal matrix both the components of  $\mathbf{x}'$  and the components of  $\xi'$  are independent and the same happens for the nonzero components of  $\xi = (\xi')^{\kappa \uparrow}$  causing  $\mathcal{E}$ , and thus also  $\mathcal{X}$  to be diagonal. In this case,  $\mathbf{x}$  will only be  $\kappa$ -sparse but not localized. In fact, by recalling (1.5) we have

$$\mathcal{L}_x = \frac{\text{tr}(\mathcal{X}^2)}{\text{tr}^2(\mathcal{X})} - \frac{1}{n} = 0$$

since for diagonal correlations we have  $\text{tr}(\mathcal{X}^2) = n\mathcal{X}_{0,0}^2$  while  $\text{tr}^2(\mathcal{X}) = n^2\mathcal{X}_{0,0}^2$ .

Localization can be imposed by choosing a non-diagonal  $\mathcal{X}'$  whose features will approximately be translated into those of  $\mathcal{X}$ . In fact, since the  $\kappa$  largest components of  $\xi'$  are carried over to  $\xi$ ,  $\xi$  is the best possible  $\kappa$ -sparse approximation of  $\xi'$  and the same relationship holds between  $x$  and  $x'$ . Hence, the larger the  $\kappa$ , the more similar the behavior of  $x$  to that of  $x'$ .

Though, the relationship between the localization of  $\mathcal{X}'$  and that of  $\mathcal{X}$  is difficult to model analytically we may provide some numerical evidence on the effectiveness of this method within the specific framework that we will use in our examples.

In particular we will consider  $\mathcal{X}'$  such that  $\mathcal{X}'_{j,k} = \omega^{|j-k|}$  for some  $-1 < \omega < 1$ . As noted in Chap. 1, this means that  $x'$  is a chunk of a stationary stochastic process with power spectrum

$$\Psi(f) = \frac{1 - \omega^2}{1 + \omega^2 - 2\omega \cos(2\pi f)}$$

that assumes a high-pass profile for  $-1 < \omega < 0$ , a flat/white profile for  $\omega = 0$ , and a low-pass profile for  $0 < \omega < 1$ . With some calculations one gets

$$\mathcal{L}_{x'} = \frac{2}{n^2} \sum_{j=1}^{n-1} j\omega^{2(n-j)} = \frac{2\omega^2}{n} \frac{n(1 - \omega^2) + \omega^{2n} - 1}{n(1 - \omega^2)^2} \quad (2.5)$$

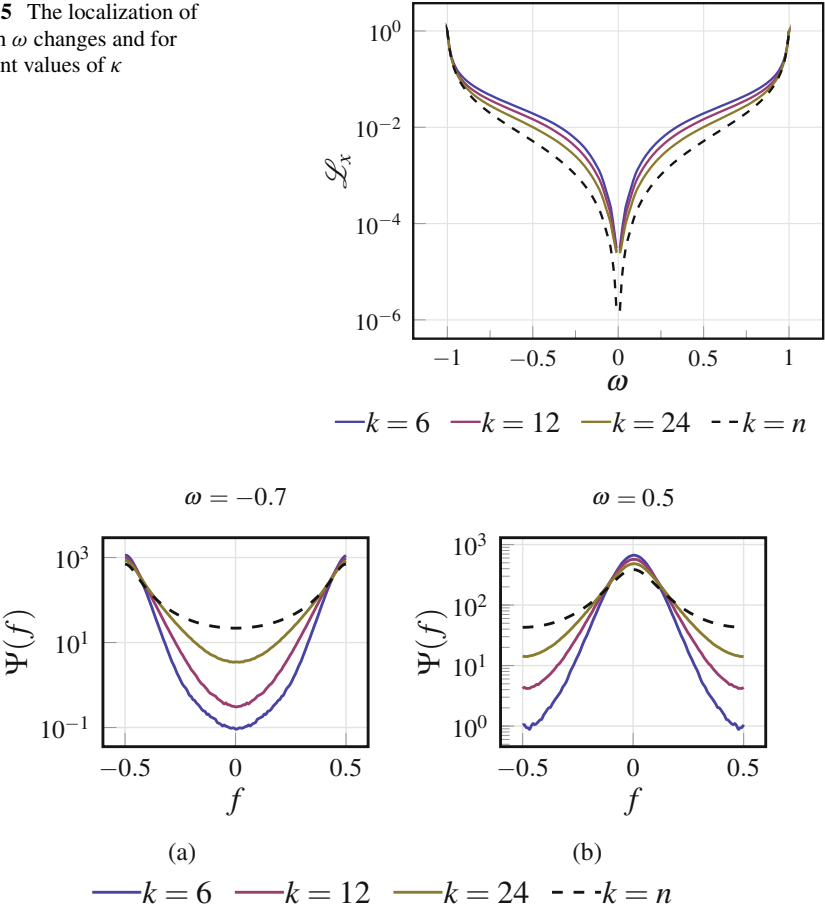
Assume now  $n = 128$ , and  $D$  as the orthonormal Discrete Cosine Transform (DCT) basis. By generating a large amount of sample vectors  $x'$  and thus  $x$  we may estimate their localization and the power spectrum of the process from which they are taken. The result of such estimations is reported in Figs. 2.5 and 2.6 for different values of  $\omega$  and sparsity  $\kappa$  (remember that  $\kappa = n$  implies  $x = x'$ , i.e.,  $x$  is a Gaussian random vector with an exponential correlation controlled by the decay  $\omega$ ).

In particular Fig. 2.5 shows how  $\mathcal{L}_x$  changes when  $\omega$  changes. Note that increasing  $|\omega|$  increases the localization of the generated signal  $x$ . As far as *what kind* of localization is conferred to  $x$ , Fig. 2.6 shows that when  $x'$  is low-low pass also  $x$  is low-pass, and vice versa.

Overall, our generation methods prove itself to be a practical way of ensuring sparsity while at least qualitatively controlling the localization of the signal and, as such, will be used in all the non-real-world examples of this volume.

To keep such examples not too far from realistic conditions, we refer to Table 1.1 and focus on processes with localizations compatible with those of real-world signals. This helps defining some prototype signals that are reported in Table 2.6.

**Fig. 2.5** The localization of  $\mathbf{x}$  when  $\omega$  changes and for different values of  $\kappa$



**Fig. 2.6** The spectrum of  $\mathbf{x}$  in a low-pass (b) and high-pass (a) case for different values of  $\kappa$

**Table 2.6** Definition of prototype signals used in the toy examples

Signal name	$\mathcal{L}_x$	$\kappa$	$\omega$
ZL: $\mathcal{L}_x = 0$ —white	0	6	0
		12	
		24	
LL: low $\mathcal{L}_x$	0.02	6	$\pm 0.509$
		12	$\pm 0.584$
		24	$\pm 0.669$
ML: medium $\mathcal{L}_x$	0.06	6	$\pm 0.810$
		12	$\pm 0.853$
		24	$\pm 0.878$
HL: high $\mathcal{L}_x$	0.2	6	$\pm 0.959$
		12	$\pm 0.964$
		24	$\pm 0.966$

Since most of our discussion hinges on the design of the matrix  $\mathbf{A}$  producing the compressed measurements  $\mathbf{y} = \mathbf{A}\mathbf{x}$ , all the examples will address a specific design option or compare a number of them.

To do so we will rely on signals  $\mathbf{x}$  generated as above and simulate the acquisition process by first perturbing them with a random vector  $\boldsymbol{\eta}^x$  made of independent, zero-mean Gaussian components whose variance is adjusted to match a prescribed Intrinsic Signal-to-Noise Ratio

$$\text{ISNR[dB]} = 20 \log_{10} \left( \frac{\|\mathbf{x}\|_2}{\|\boldsymbol{\eta}^x\|_2} \right)$$

Such a perturbation is injected to simulate inaccuracies in the acquisition stages, including the possible quantization.

The perturbed signal is then used to produce measurements by using the matrix  $\mathbf{A}$  under assessment and obtaining  $\mathbf{y} = \mathbf{A}(\mathbf{x} + \boldsymbol{\eta}^x)$ . When not explicitly declared otherwise, we will produce the reconstructed signal to be matched against the true signal by feeding  $\mathbf{y}$ ,  $\mathbf{A}$ , and ISNR into the functions provided by the SPGL1 package mentioned in the previous section implementing either the BP or BPDN method, whose robustness and ease of use make it the ideal candidate for the concoction of examples.

## 2.4 Practical Performance

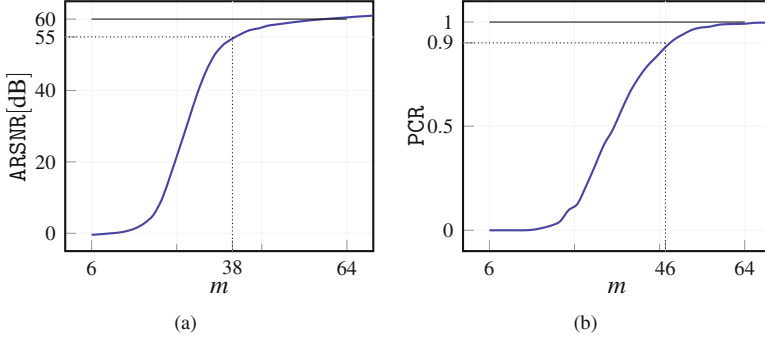
The first section of this chapter shows that, when not modeled from a worst-case point of view, CS is a promising technique that may allow to reconstruct an  $n$ -dimensional signal  $\mathbf{x}$  from  $m$  scalar measurements in a vector  $\mathbf{y}$  with  $m \ll n$ .

To give a quantitative appreciation of what can be achieved, assume that  $\mathbf{x}$  is  $n$ -dimensional with  $n = 128$ , that is  $\kappa = 6$  sparse with respect to a DCT orthonormal basis and that is generated as described before with  $\omega = 0$  and  $\text{ISNR[dB]} = 60$  dB.

Take  $\mathbf{A} \sim \text{RGE (iid)}$  and, for each value of  $m$  from 6 to 64 perform a Montecarlo simulation. For each trial compute the RSNR to accumulate a profile of ARSNR as a function of  $m$ . By fixing  $\text{RSNR[dB]}_{\min} = \text{ISNR[dB]} - 5 \text{ dB} = 55 \text{ dB}$ , we may also estimate the PCR for each  $m$ . The result is reported in Fig. 2.7.

Since both ARSNR and PCR are the-larger-the-better merit figures, the sigmoidal trends in both plots are the practical implication of theorems like Theorem 2.4. In fact, coherently with what theory says, there is some critical value of  $m$  after which performance dramatically increases, giving rise to what is often called *phase transition*.

Beyond reflecting theoretical results, plots like those in Fig. 2.7 give a quantitative appreciation of achievable *compression*. For example, Fig. 2.7a shows that for  $m = 64$  the ARSNR slightly exceeds the  $\text{ISNR} = 60$  dB and thus indicates that the original signal, whose dimensionality is  $n = 128$ , can be acquired with a compression ratio  $\text{CR} \simeq 2$  with no loss of accuracy (actually with a little amount



**Fig. 2.7** Montecarlo assessment of performance for a classical CS system: when the number of measurements increases both the ARSNR **(a)** and PCR **(b)** increase

of denoising). Yet, one may decide that an  $\text{ARSNR} = 55 \text{ dB}$  is enough for the application at hand and derive from the same plot that  $m = 38$  measurements are enough to meet the specification, increasing the compression ratio to  $\text{CR} \simeq 3.4$ .

Clearly, this concerns average performance. A stricter point of view would be to require that  $\text{RSNR} = 55 \text{ dB}$  is not achieved on average but at least 90% of the times. Since Fig. 2.7b estimates the probability that  $\text{RSNR}$  exceeds that threshold as a function of  $m$ , one gets that this more stringent specification can be met sizing the system with  $m = 46$ , that still gives  $\text{CR} \simeq 2.8$ .

This is a somehow impressive performance since it places well apart from the worst-case scenarios that were addressed in deriving the guarantees. As an example, Theorem 1.7, specialized to the case in which  $\mathbf{x}$  is perfectly  $\kappa$ -sparse with respect to an orthonormal basis, says that the error between the original signal  $\mathbf{x}$  and its reconstruction  $\hat{\mathbf{x}}$  can be bounded as

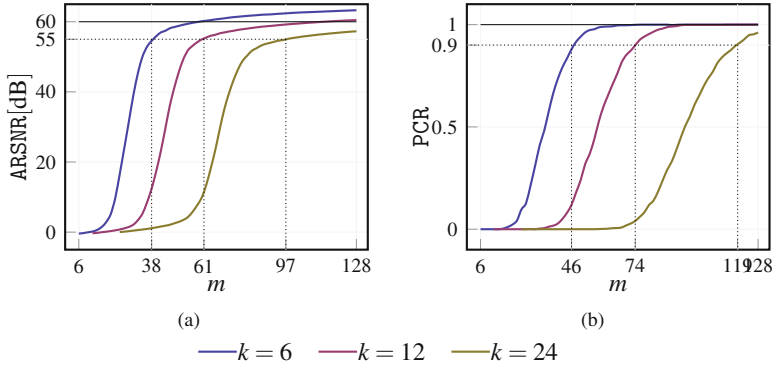
$$\|\hat{\mathbf{x}} - \mathbf{x}\|_2 = \|\mathbf{D}\hat{\boldsymbol{\xi}} - \mathbf{D}\boldsymbol{\xi}\|_2 = \|\hat{\boldsymbol{\xi}} - \boldsymbol{\xi}\|_2 \leq 4 \frac{\sqrt{1 + \delta_{2k}}}{1 - (\sqrt{2} + 1)\delta_{2k}} \epsilon$$

where  $\epsilon$  is such that  $\|\boldsymbol{\eta}\|_2 \leq \epsilon$ , and  $\delta_{2k}$  is the RIC of  $\mathbf{A}$ . For  $\delta_{2k} \geq 0$ , the coefficient of  $\epsilon$  is monotonically increasing and thus, even in the best possible conditions, the guarantee of Theorem 1.7 on the  $\text{RSNR}$  is

$$\text{RSNR}[\text{dB}] = 20 \log_{10} \left( \frac{\|\mathbf{x}\|_2}{\|\hat{\mathbf{x}} - \mathbf{x}\|_2} \right) \geq 20 \log_{10} \left( \frac{\|\mathbf{x}\|_2}{4 \|\boldsymbol{\eta}\|_2} \right) \geq \text{ISNR}[\text{dB}] - 12 \text{ dB}$$

that, due to its worst-case nature, gives little hint on the fact that, for example, a small average denoising effect can be obtained.

Figure 2.8 shows the trends of same merit figures when the signal to acquire is either  $\kappa = 6$ -sparse, or  $\kappa = 12$ -sparse, or  $\kappa = 24$ -sparse. Clearly, since  $\kappa$  is the minimum number of scalars that are needed to identify  $\mathbf{x}$ , a progressively larger



**Fig. 2.8** Montecarlo assessment of performance for a classical CS system: when the number of measurements increases both the ARSNR (a) and PCR (b) increase, though with trends depending on the sparsity  $\kappa$

**Table 2.7** Numerical matching between the asymptotic trend in (1.29) and the empirical evidence of Fig. 2.8. The increase in the sparsity of the signal  $\kappa$  implies an increase in the minimum number  $m^*$  of measurements needed to achieve a certain performance that is compared with the  $O(\kappa \log(n/\kappa))$  trend

$\kappa$	ARSNR $\geq 55$ dB		PCR $\geq 0.9$	
	$m^*$	$\frac{m^*}{\kappa \log_2(n/\kappa)}$	$m^*$	$\frac{m^*}{\kappa \log_2(n/\kappa)}$
6	38	1.43	46	1.74
12	61	1.49	74	1.81
24	97	1.67	119	2.05

number of measurement is needed to achieve a good signal reconstruction and the corresponding curves move to the right.

As an example, to obtain  $\text{ARSNR} \geq 55$  dB one needs at least  $m^* = 38$  measurements when the signal is  $\kappa = 6$ -sparse, but  $m^* = 61$  measurements if the signal to reconstruct is  $\kappa = 12$ -sparse, and  $m^* = 97$  measurements for  $\kappa = 24$ -sparse signals. Though the trend of  $m^*$  against  $n$  and  $\kappa$  is identified only in asymptotic terms by (1.28) and (1.29), it may be used as a rough estimate of  $m^*$  even in finite cases.

In fact, by looking at Table 2.7 one is tempted to adopt as a first sizing criterion  $m^* = c\kappa \log_2(n/\kappa)$  with a constant  $c$  in the range  $2 \leq c \leq 3$ .

Though all this may seem a success, from an engineering point of view it is only a starting point. In fact, performances like those in Fig. 2.7 are estimated for a system in which  $\mathbf{A}$

- has entries that are infinite precision and unbounded;
- is maximally random within the variance constraint on its entries, and thus is completely agnostic both of its role and of its optimization possibilities.

Yet, any real-world implementation of the multiplication of  $\mathbf{x}$  by  $\mathbf{A}$  will imply a finite-range calculation with a limited precision, either because of noise if the implementation is analog, or because of quantization if the implementation is digital.

Moreover, instead of simply *accepting* measurements as they happen to be computed by a maximally random policy, one may try to look for measurements that best identify the signal itself so to squeeze as much information as possible in the  $m < n$  scalars that will represent  $\mathbf{x}$ . The hope of this quest for *good* measurements is that a smaller number of substantial pieces of information can do the same job of a larger number of purely random looks at the signal.

Leaving the finite-range/finite-precision issue to a following chapter, note that this second aim seems to go against a quite commonly accepted idea, suggestively indicated as *democracy*, that each measurement carries roughly the same amount of information about the signal being acquired. The mathematical foundation of this idea is solid, it has to do with the RIC of the matrices  $\mathbf{A}$  and with how such constant changes when few rows are dropped: it turns out that when the number of surviving rows is still larger than the minimum number needed to guarantee signal reconstruction, then which row was discarded has little effect on the RIC constant of the resulting matrix.

Yet, the practical effects of such a formal development are negligible for at least two reasons. The first is that we have seen how RIC-based performance bounds are to be taken only as guarantees since they are so loosely correlated with real-world performance that their use as a design criterion is ineffective. In this case, a small change in the RIC of  $\mathbf{A}$  implies a small change in the performance guarantee but says nothing on the change of the actual performance.

The second is that when seeking for optimally designed CS stages, one never moves from the minimum  $m$  needed for a correct reconstruction far enough to be able to speculate about dropping some measurements and still working above that limit. The aim of system optimization is to push  $m$  as low as possible.

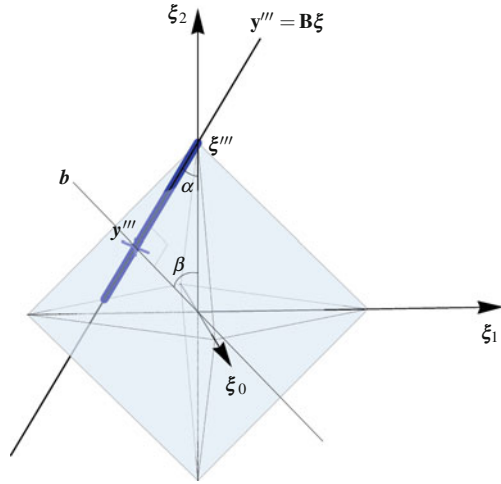
## 2.5 Countering the Myth of Democracy and Paving the Way for Practical Optimization

Beyond the considerations in the previous section, the fact that measurement *democracy* is a myth incorrectly inferred from a sound mathematical result can be demonstrated with an easy formal argument if we go back to the simplified setting that characterizes the first sections of this chapter: no noise (i.e.,  $\text{ISNR} = \infty$ ), and straightforward BP for reconstruction.

Within such a framework, we benefit from a powerful geometric insight on the reasons why the basic reconstruction strategy is successful and when it fails. In particular we may concentrate on failures and have a second look at a slightly rotated and simplified version of Fig. 2.3, that is Fig. 2.9.



**Fig. 2.9** A rotated and simplified version of Fig. 2.3 that highlights the relative positions of the signal  $\xi''$ , the projection  $y'''$  and the face of  $S_3^1(\|\xi'''\|_1)$



In that figure, it is easy to verify that the solution to BP is not unique (all the points on the thick blue segment are possible reconstructions of the original signal) due to the fact that the projection plane contains a direction  $\mathbf{b}$  that is orthogonal to one of the 2-dimensional facets of  $S_3^1(\|\xi'''\|_1)$  to which  $\xi'''$  itself belongs.

This may be reworded saying that one of the vectors onto which the signal is projected (one of the rows  $\mathbf{b}$  of the matrix  $\mathbf{B} = \mathbf{AD}$ ) forms with the signal an angle  $\beta$  such that its complement  $\alpha = \pi/2 - \beta$  is equal to the angle between the signal  $\xi$  and a facet of  $S_3^1(\|\xi'''\|_1)$ .

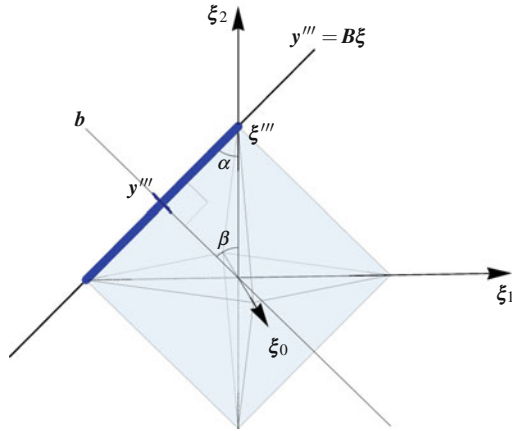
In this case  $\alpha = \arccos(\sqrt{2/3})$  and if we might ensure that the rows of  $\mathbf{B}$  avoid forming with the signal an angle  $\beta = \pi/2 - \arccos(\sqrt{2/3})$ , a case like the one depicted in Fig. 2.9 never occurs and signals like  $\xi'''$  are correctly reconstructed.

Clearly, other *bad* cases may happen. As an example, Fig. 2.10 shows that there is another choice of  $\mathbf{b}$  that prevents BP from retrieving the original signal. Again, the reason is that the angle  $\alpha$  between the signal and a facet of  $S_3^1(\|\xi'''\|_1)$  (in this case it is a 1-dimensional facet) is complementary to the one  $\beta$  between the signal and the direction  $\mathbf{b}$  along which we are projecting. In this case  $\alpha = \beta = \arccos(\sqrt{1/2}) = \pi/4$ .

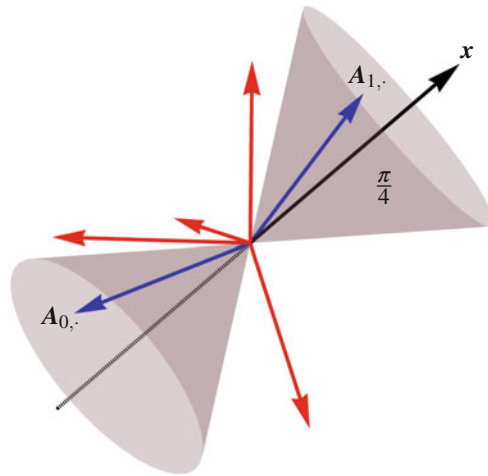
If we avoided choosing directions  $\mathbf{b}$  whose angle the signal angles is one of the two computed above, both *bad* cases would be prevented.

Though the detailed proof is out of the scope of this volume, in the general  $n$ -dimensional case, when the signal is  $\kappa$ -sparse, angles between  $\beta^{\min} = \pi/2 - \arccos(1/\sqrt{1+\kappa})$  and  $\beta^{\max} = \pi/2 - \arccos(\sqrt{\frac{n-\kappa}{n-\kappa+1}})$  must be avoided. In our case  $n = 3$  and  $\kappa = 1$  so that  $\beta^{\min}$  and  $\beta^{\max}$  boil down to  $\pi/4$  and  $\pi/2 - \arccos(\sqrt{2/3})$  computed before. To be on the safe side and not to take too subtle decisions depending on  $n$  and  $\kappa$ , all angles in  $[\pi/4, \pi/2]$  should be avoided.

**Fig. 2.10** Another *bad* choice of a direction  $\mathbf{b}$  to use for projection. Also in this case BP cannot reconstruct the original signal as all the points in the *thick blue segment* are possible solutions

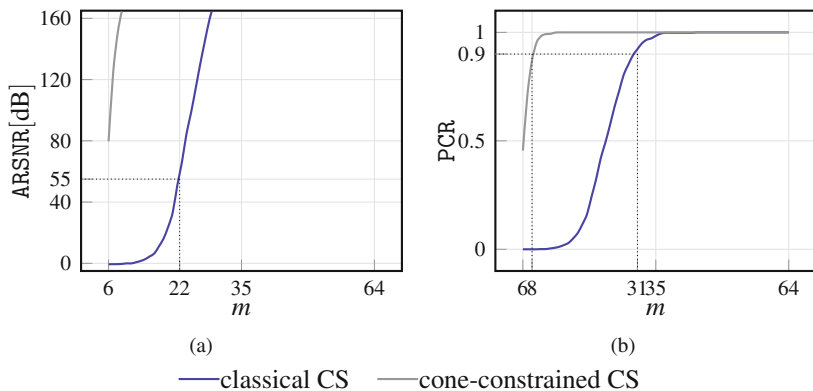


**Fig. 2.11** Among many candidate vectors randomly pointing in space, only the two falling in the cone whose axis is  $\mathbf{x}$  and whose aperture is  $\pi/4$  become the rows  $\mathbf{A}_{0,\cdot}$  and  $\mathbf{A}_{1,\cdot}$  of the matrix  $\mathbf{A}$  used for acquisitions



Hence, to ensure maximum reconstruction performance in a noiseless environment, it is advisable to select matrices  $\mathbf{B}$  whose rows form with  $\xi$  an angle strictly smaller than  $\pi/4$ . When  $\mathbf{D}$  is an orthonormal basis, this directly translates into a prescription for the angle between the rows of  $\mathbf{A} = \mathbf{B}\mathbf{D}^\top$  and the signal  $\mathbf{x} = \mathbf{D}\xi$ . Such a prescription translates into a simple geometric criterion: rows may be generated as  $n$ -dimensional vectors whose entries are independent random variables  $\sim \mathcal{N}(0, 1)$  but are included in  $\mathbf{A}$  only if their angle with  $\mathbf{x}$  is less than  $\pi/4$ . Such a method will be indicated as *cone-constrained CS* as accepted rows are vectors falling in the cone whose axis is  $\mathbf{x}$  and whose aperture is  $\pi/4$  as exemplified in Fig. 2.11.

To have a practical appreciation of how much this criterion affects performance we may adopt the same simulation setting as above in the noiseless  $\text{ISNR} = \infty$  case and substituting BPDN with BP implemented as a purely linear optimization



**Fig. 2.12** Montecarlo assessment of performance for classical CS (blue) and cone-constrained CS (gray): the ideal cone-constrained CS clearly exhibits far better performance. (a) performance in terms of ARSNR while (b) is for PCR

problem (1.32). In these conditions we simulate the performance of classical CS, and cone-constrained CS. The results are shown in Fig. 2.12.

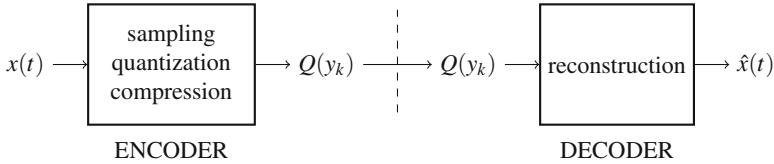
The absence of noise clearly improves reconstruction performance of classical CS. By comparing Fig. 2.7 with Fig. 2.12 we get that an average quality ARSNR = 55 dB can be reached with  $m = 22$  measurements instead of  $m = 38$  measurements, and an RSNR = 55 dB can be guaranteed 90% of the times with  $m = 31$  measurements instead of  $m = 48$  measurements.

Yet, cone-constrained CS has definitely better performance since the average reconstruction quality never falls below 80 dB and RSNR = 55 dB can be guaranteed 90% of the times with only  $m = 8$  measurements.

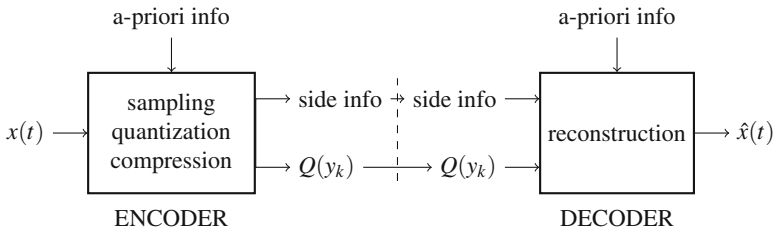
Overall, the measurements we select are clearly carrying more information about the signal with respect to measurements picked randomly, and no democracy exists in the real-world. Although this may be taken as a discomfoting truth from a social point of view, it is actually extremely good news from the point of view of engineering of CS. In fact, when not all the options are equally good, optimization may be called into play to look for the best design alternatives.

Regrettably, cone-constrained CS is only a theoretical tool since it has no concrete chance to be implemented. To understand why, we have to spend a few words on some very high-level implementation constraints that are at the base of any successful application of CS.

Though it is true that this book focuses on the design of CS stages according to the scheme of Fig. 1.2 we cannot avoid to place such an acquisition subsystem in a slightly more general perspective. This is what Fig. 2.13 does considering the same quantities as in Fig. 1.2. All the acquisition stages (sampling, quantization, compression) can be seen as a single block that *encodes* the analog waveform into the subsufficient-rate sequence of digital scalars  $Q(y_k)$ . Such sequence is passed to some other subsystem that is interested in knowing  $x(t)$  and *decodes* the sequence  $Q(y_k)$  into an approximation  $\hat{x}(t)$ .



**Fig. 2.13** A higher-level view of the role of signal acquisition



**Fig. 2.14** The encoder-channel-decoder view with additional signal paths

The higher-level point of view reveals an encoder–decoder structure that highlights the fact that the only continuous communication between the encoder and the decoder is the subsufficient sequence, i.e., in principle, no other signal dependent information is passed from the acquisition subsystem to the subsystem that uses the acquired signal.

In terms of a CS acquisition mechanism, this means that, for example, the rows of  $\mathbf{A}$  are not communicated to the decoder, that must be able to know them independently. This is why, a more realistic view at the scheme in Fig. 2.13 should comprise few other details.

First, the encoder and the decoder sides must share some a priori information. If, for example,  $\mathbf{A}$  is fixed, then it must enter the design of both sides. Alternatively, if  $\mathbf{A}$  is a time varying instance of a random matrix ensemble (as in our examples), the encoder and the decoder may share the design of a reproducible pseudorandom number generator and the initial state from which it works. In this case the operations of encoder and decoder must be synchronized thus implying a small amount of side information to be transferred from encoder to decoder further to the subsufficient sequence  $Q(y_k)$ . The resulting more realistic view of the acquisition system is given in Fig. 2.14. Clearly, for the compression scheme to be effective, the total transferred information (the subsufficient sequence plus the side information) must amount to less bits than what would be needed by the sheer transmission of a sufficient sequence of samples.

This is the main reason why cone-constrained CS cannot be effectively employed. In fact, what we may do to apply the method in practice is to deploy two identical copies of a pseudorandom number generator both at the encoder and

the decoder, synchronize them and let them run to produce candidate rows for the matrix  $\mathbf{A}$ . The encoder tests each of them and accepts only the first  $m$  of them whose angle with  $\mathbf{x}$  is less than  $\pi/4$  to build  $\mathbf{A}$ . Then, it computes  $\mathbf{y} = \mathbf{A}\mathbf{x}$  and communicates to the decoder both the vector  $\mathbf{y}$  and the side information needed to identify the rows it used.

If we assume that to find  $m$  rows one must examine  $M$  candidates, the number of bits of side information is  $\left\lceil \log_2 \binom{M}{m} \right\rceil$  since our task is to identify a specific subset of  $m$  elements among  $M$  possible candidates. Overall, the amount of information that must be transferred from the encoder to the decoder is  $m\mathbf{b}_y + \left\lceil \log_2 \binom{M}{m} \right\rceil$ , where  $\mathbf{b}_y$  is the number of bits used for each sample of the subsufficient sequence  $Q(y_k)$ . This must be compared with the straightforward option of quantizing each samples with  $\mathbf{b}_x$  bits so that the bitwise compression ratio is

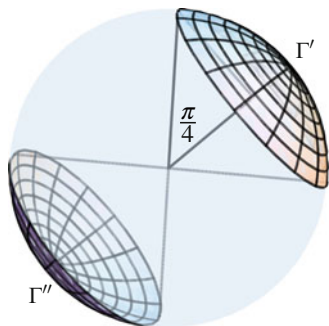
$$\text{CR}^{\text{bit}} = \frac{n\mathbf{b}_x}{m\mathbf{b}_y + \left\lceil \log_2 \binom{M}{m} \right\rceil} \quad (2.6)$$

Regrettably, the ratio between  $m$  and  $M$  suffers from a well-known effect of dimensionality on the shape of  $S_n^2$  spheres. Assuming that the candidate rows span all the possible angles uniformly (this is what happens, for example, if their entries are independent normals), the probability that one of them falls within the proper cone is equal to the ratio between the measure of the surface of the two spherical caps  $\Gamma' \cup \Gamma''$  illustrated in Fig. 2.15 and the measure of the surface  $\partial S_n^2$  of the whole sphere  $S_n^2$ .

From [15] we get that such a ratio is

$$\frac{\mu(\Gamma' \cup \Gamma'')}{\mu(\partial S_n^2)} = B_{\sin^2(\pi/4)} \left( \frac{n-1}{2}, \frac{1}{2} \right) = B_{1/2} \left( \frac{n-1}{2}, \frac{1}{2} \right)$$

**Fig. 2.15** The spherical caps whose surface is proportional to the probability of generating a random measurement falling into the  $\pi/4$  cone



that uses the incomplete regularized beta function

$$B_{\zeta}(p, q) = \frac{\int_0^{\zeta} t^{p-1} (1-t)^{q-1} dt}{\int_0^1 t^{p-1} (1-t)^{q-1} dt}$$

from which we derive that  $B_{1/2}(\frac{n-1}{2}, \frac{1}{2}) \leq 2^{-\frac{n-1}{2}}$  for  $n \geq 1$  is decreasing not less than exponentially with  $n$ . This means, for example, that the probability of a candidate 128-dimensional row of falling into the  $\pi/4$  cone whose axis is any given signal  $\mathbf{x}$  is less than  $7.6 \times 10^{-21}$ .

Assume now that we want to guarantee that  $\text{RSNR} \geq 55$  dB at least 90% of the times. From Fig. 2.12 we get that  $m = 8$  measurements are enough. Yet, the average number of independent candidate rows to evaluate before accumulating  $m = 8$  measurements is  $8/(7.6 \times 10^{-21}) = 1.1 \times 10^{21}$ , and the side information that needs to be communicated amounts to 544 bit. If we assume  $b_y = 12$  (a sensible choice to achieve  $\text{RSNR} = 55$  dB), the total number of bits needed to encode the  $n = 128$ -dimensional window is  $544 + 8 \times 12 = 640$  bit.

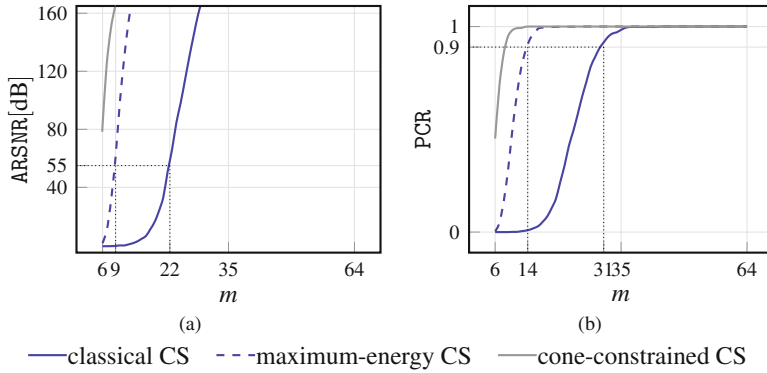
Without compression, we may roughly estimate the  $\text{RSNR}$  achieved by the straightforward quantization of each sample assuming that signal behaves almost sinusoidally so that  $\text{RSNR} = 6.02b_x + 1.76$  dB where  $b_x$  is the number of bits used for each sample. To have  $\text{RSNR} = 55$  dB we may set  $b_x = 9$  so that the total number of bits would be  $128 \times 9 = 1152$  bit. Hence, the bitwise compression ratio (2.6) of cone-constrained CS is  $\text{CR}^{\text{bit}} = 1152/640 \simeq 1.8$ .

Note that, if we decided to use every row produced by the generator we would not need to send any side information beyond an initial synchronization, while Fig. 2.12 tells us that the same performance level as before would be guaranteed by  $m = 21$  measurements, for a total of  $mb_y = 31 \times 12 = 372$  bit and a corresponding bitwise compression ratio  $\text{CR}^{\text{bit}} = 1152/372 \simeq 3.1$ .

All this said, there is no point in trying an implementation of cone-constrained CS since it does not give any real advantage with respect to purely random CS which also enjoys a much smaller computational burden (even if generating and testing a candidate row took a single nanosecond, accumulating 8 measurements in the  $\pi/4$  cone would take more than 33000 years!<sup>1</sup>).

However, what we are left with is an intuition that a possible criterion to increase the amount of information that a measurement  $y$  carries about the signal  $\mathbf{x}$  is to obtain it as  $y = \mathbf{a}^T \mathbf{x}$  using a vector  $\mathbf{a}$  lying on straight line whose angle with the straight line containing  $\mathbf{x}$  is ‘small’, whatever this may signify. From here on what we do can only be intuitively justified but, as we will see in the more applicative

<sup>1</sup>To avoid this curse of dimensionality, the simulations leading to Fig. 2.12 had to generate rows of  $\mathbf{A}$  by properly modulating the length of random rotations of  $\mathbf{x}$  itself with angle smaller than  $\pi/4$ , the conceptual sieving procedure being totally unfeasible.



**Fig. 2.16** Montecarlo assessment of performance for classical CS (solid blue), maximum-energy CS (dashed blue), and cone-constrained CS (solid gray): maximum-energy CS is not as performing as cone-constrained CS but it outperforms classical CS. (a) performance in terms of ARSNR while (b) is for PCR

chapters of this volume, gives raise to a powerful heuristic criterion supporting a well-defined and effective design flow for practical CS acquisition.

The first trivial remark is that, given two non-collinear vectors  $\mathbf{v}'$  and  $\mathbf{v}''$ , forming an angle  $\widehat{\mathbf{v}'\mathbf{v}''}$ , the angle between the straight lines containing them is  $\min\{\widehat{\mathbf{v}'\mathbf{v}''}, \pi - \widehat{\mathbf{v}'\mathbf{v}''}\}$ . Hence, such an angle gets smaller when  $\widehat{\mathbf{v}'\mathbf{v}''}$  either goes to 0 or  $\pi$ , i.e., when the absolute values  $\cos^2(\widehat{\mathbf{v}'\mathbf{v}''})$  increases. Since  $y^2 = (\mathbf{a}^\top \mathbf{x})^2 = \|\mathbf{a}\|_2^2 \|\mathbf{x}\|_2^2 \cos^2(\widehat{\mathbf{a}\mathbf{x}})$  and  $\mathbf{x}$  is assigned, if we may assume that all the rows of  $\mathbf{A}$  have approximately the same length, smaller angles correspond to higher energies of the measurement  $y$ .

A new, heuristic method is naturally born from these considerations. In a system analogous to what has been sketched for cone-constrained CS, let the row generator produce  $M$  candidates. Then compose the matrix  $\mathbf{A}$  with the rows  $\mathbf{a}$  corresponding to the  $m$  largest values of  $(\mathbf{a}^\top \mathbf{x})^2$ . The measurements are computed as  $\mathbf{y} = \mathbf{A}\mathbf{x}$  and passed to the decoder. This strategy is named *maximum-energy CS*.

In this new configuration both  $M$  and  $m$  are degrees of freedom. This gives us some control on the amount of bits spent on side information  $\left\lceil \log_2 \binom{M}{m} \right\rceil$ , an amount that must be traded with the quality of the reconstruction. In this case we do not have a theoretical background allowing to anticipate reconstruction performance and we have to rely on simulations. If we do so, we may add a track to Fig. 2.12 and obtain Fig. 2.16.

Maximum-energy CS is simulated generating  $M = 512$  candidates and taking the  $m$  largest energy measurements for  $m = \kappa = 6$  to  $m = n/2 = 64$ . Since it is only a heuristic approximation of the cone-constrained policy, performance decreases but is still much higher than that of classical CS.

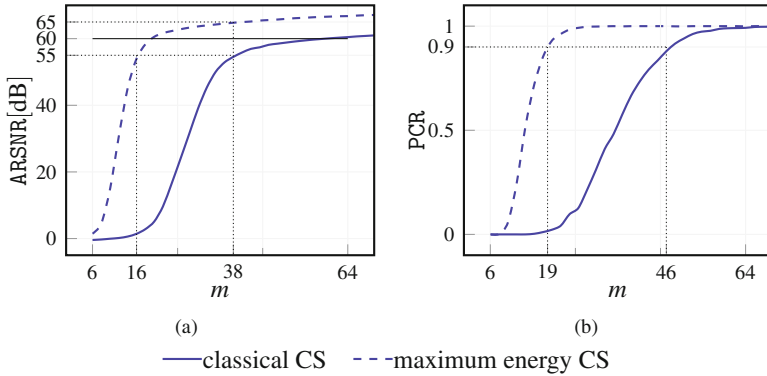
If we assume that our target reconstruction quality is  $\text{ARSNR} \geq 55$  dB, then classical CS achieves it with  $m = 22$  while maximum-energy CS requires only  $m = 9$ . These numbers allow to compute the bitwise compression ratio (2.6) for the same case as above in which  $b_x = 9$  and  $b_y = 12$ . Classical CS yields  $\text{CR}^{\text{bit}} = 1152 / (22 \times 12) = 1152 / 264 \simeq 4.4$ . Maximum-energy CS yields  $\text{CR}^{\text{bit}} = 1152 / \left( 9 \times 12 + \left\lceil \log_2 \binom{512}{9} \right\rceil \right) = 1152 / 171 \simeq 6.7$  and is therefore able to provide a gain further to the mere reduction of the number of scalar measurements.

If we assume that our target reconstruction quality is to guarantee that 90% of the times we have  $\text{RSNR} \geq 55$  dB, then classical CS achieves it with  $m = 31$  while maximum-energy CS requires only  $m = 14$ . The corresponding bitwise compression ratios are  $\text{CR}^{\text{bit}} \simeq 3.1$  for classical CS and  $\text{CR}^{\text{bit}} \simeq 4.5$  for maximum-energy CS.

Overall, maximum-energy CS seems to be a good candidate to leverage the intuitive criterion we have developed for measurement quality while keeping adaptivity to a level that can be managed by adding a reasonable amount of side information.

This is true even in a noisy environment. In fact, we may go back to our original setting in which  $\text{ISNR} = 60$  dB, and keep the same configuration of maximum-energy CS to obtain curves like the ones in Fig. 2.17.

Since noise is back, performance deteriorates also for maximum-energy CS. Yet, the new method is still able to yield better bitwise compression ratios. In fact, looking at Fig. 2.17a we get that the reference quality level  $\text{ARSNR} = 55$  dB is achieved with  $m = 38$  measurements by classical CS, and with  $m = 16$  measurements by maximum-energy CS. The usual computation of  $\text{CR}^{\text{bit}}$  with  $b_x = 9$  and  $b_y = 12$  gives  $\text{CR}^{\text{bit}} = 1152 / 456 \simeq 2.5$  for classical CS and  $\text{CR}^{\text{bit}} = 1152 / 292 \simeq 3.9$  for maximum-energy CS.



**Fig. 2.17** Montecarlo comparison between performance of classical CS (*solid*) and maximum-energy CS (*dashed*): both the ARSNR (a) and PCR (b) curves are dramatically improved



Figure 2.17a also shows that the adapted method is able to provide noteworthy average denoising. For example, when classical CS achieves our reference performance level  $\text{ARSNR} = 55 \text{ dB}$ , maximum-energy CS is able to yield  $\text{ARSNR} = 65 \text{ dB} > \text{ISNR}$ . Clearly, this comes at some expense since maximum-energy CS accumulates and communicates  $\left\lceil \log_2 \binom{512}{38} \right\rceil = 192 \text{ bit}$  in addition to the 456 bit used to encode the measurements. Yet, in this case reconstruction provides an accuracy that would not be attained by simply encoding the samples.

To confirm that what we have observed so far can be of use, we may explore different types of signals with different sparsities. As before, no formal analysis is available but an extensive Montecarlo assessment can be pursued using different values of  $\kappa$  and different values of  $\omega$  in the exponential correlation signal mode we defined in Sect. 2.3.

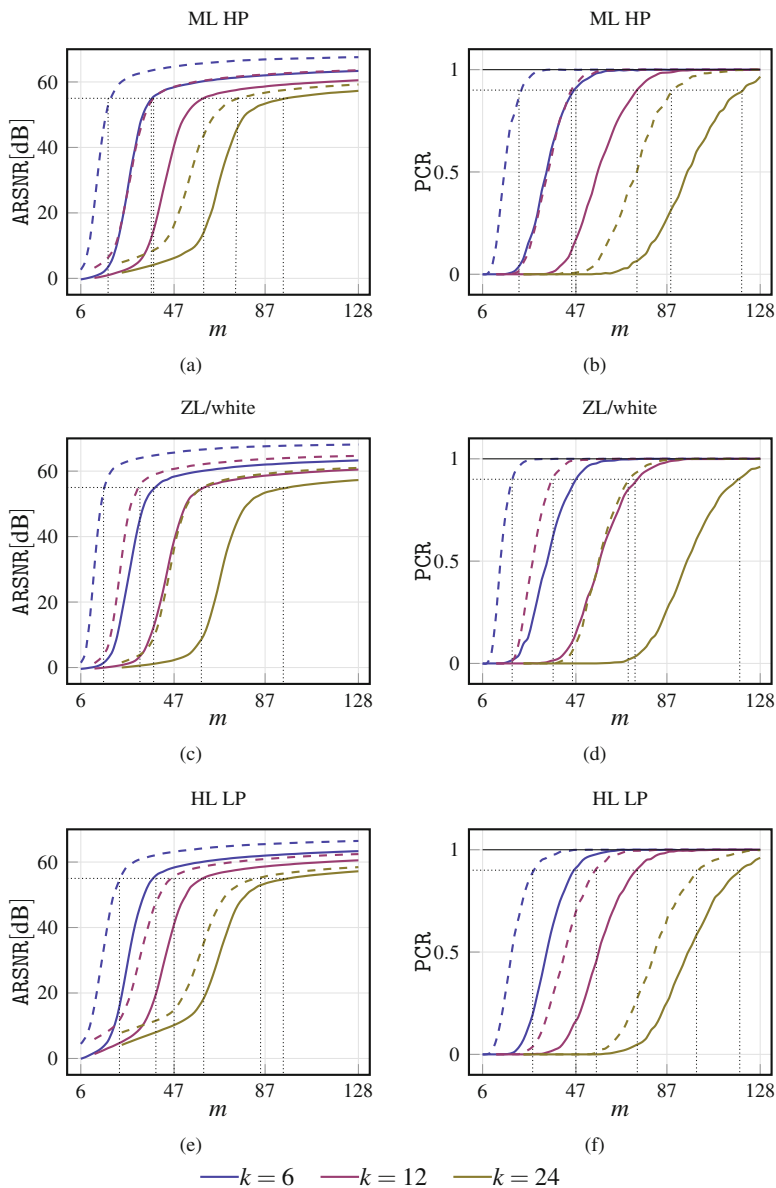
A sample of the results of such an assessment is reported in Fig. 2.18 where we give the trends of both ARSNR and PCR (with target  $\text{RSNR} = 55 \text{ dB}$ ) when  $n = 128$  and  $\kappa \in \{6, 12, 24\}$  while trying a high-pass signal with medium localization (ML HP), a white signal, and a low-pass signal with large localization (HL LP) as defined in Table 2.6. The white case with  $\kappa = 6$  is the same as the one in Fig. 2.17.

In terms of the minimum number of measurements needed to match a certain reconstruction quality, the improvement of maximum-energy CS over classical CS is undoubtable. The fact that it may result in a better bitwise compression ratio must be checked case by case. From Fig. 2.18 it is clear that the performance of classical CS is independent of localization while that of maximum-energy CS is not. Hence, we may focus on the white case and summarize our quantifications in Table 2.8.

Overall, the maximum-energy criterion seems to behave rather well. Yet, its implementation has two drawbacks that may limit practical application. First, one has to compute many more measurements ( $M = 512$  in our examples) than what is then communicated to the decoder. If the cost of computing a measurement is not negligible, this may have an impact on the encoder complexity. Due to the dimensionality effect that we already noted when analyzing cone-constrained CS, such an impact is expected to dramatically increase as  $n$  increases. Moreover, side information must be computed and communicated to the decoder further to measurements and this may also imply an increased encoder complexity.

Intuitively, this potentially increased complexity depends on the fact that maximum-energy CS automatically adapt itself to the specific instance  $\mathbf{x}$  of the signal it is acquiring.

In the next chapter, we will see that a slightly less performing method can be devised which, leveraging the same measurement energy principle and adapting to the class of signals to acquire rather than to a specific instance, allows to increase performance while keeping encoder complexity to a minimum.



**Fig. 2.18** Montecarlo comparison between performance of classical CS (solid) and maximum-energy CS (dashed): both the ARSNR (a) and PCR (b) curves are dramatically improved in all configurations

**Table 2.8** Bitwise compression ratios of classical CS and maximum-energy CS when  $b_x = 8$  and  $b_y = 12$  and for two different reconstruction quality requirements. Straightforward encoding of  $n = 128$  samples would require 1152 bit

ARSNR = 55 dB							
$\kappa$	Maximum-energy CS				Classical CS		
	$m$	$M$	$mb_y + \lceil \log_2 \binom{M}{m} \rceil$	$CR^{\text{bit}}$	$m$	$mb_y$	$CR^{\text{bit}}$
6	16	512	292	3.9	38	452	2.5
12	32	512	553	2.1	59	708	1.6
24	59	512	968	1.2	95	1140	1.0
PCR = 0.9							
$\kappa$	Maximum-energy CS				Classical CS		
	$m$	$M$	$mb_y + \lceil \log_2 \binom{M}{m} \rceil$	$CR^{\text{bit}}$	$m$	$mb_y$	$CR^{\text{bit}}$
6	19	512	342	3.9	47	564	2.5
12	37	512	632	1.8	73	876	1.3
24	70	512	1131	1.0	119	1428	0.81

## References

1. S. Becker, J. Bobin, E.J. Candès, NESTA: a fast and accurate first-order method for sparse recovery. *SIAM J. Imag. Sci.* **4**(1), 1–39 (2011)
2. E. van den Berg, M.P. Friedlander, Probing the Pareto frontier for basis pursuit solutions. *SIAM J. Sci. Comput.* **31**(2), 890–912 (2008)
3. T. Blumensath, M.E. Davies, Iterative thresholding for sparse approximations. *J. Fourier Anal. Appl.* **14**(5–6), 629–654 (2008)
4. D. Bortolotti et al., Energy-aware bio-signal compressed sensing reconstruction on the WBSN-gateway. *IEEE Trans. Emerg. Top. Comput.* **PP**(99), 1–1 (2016)
5. P.L. Combettes, J.-C. Pesquet, A proximal decomposition method for solving convex variational inverse problems. *Inverse Prob.* **24**(6), p. 065014 (2008)
6. I. Daubechies et al., Iteratively reweighted least squares minimization for sparse recovery. *Commun. Pure Appl. Math.* **63**(1), 1–38 (2010)
7. D.L. Donoho, Neighborly Polytopes and Sparse Solution of Underdetermined Linear Equations, Technical report, Department of Statistics, Stanford University, 2005
8. D.L. Donoho, J. Tanner, Counting faces of randomly projected polytopes when the projection radically lowers dimension. *J. Am. Math. Soc.* **22**(1), 1–53 (2009)
9. D.L. Donoho, J. Tanner, Observed universality of phase transitions in high-dimensional geometry with implications for modern data analysis and signal processing. *Philos. Trans. R. Soc. Lond. A Math. Phys. Eng. Sci.* **367**(1906), 4273–4293 (2009)
10. D.L. Donoho, J. Tanner, Precise undersampling theorems. *Proc. IEEE* **98**(6), 913–924 (2010)
11. I.F. Gorodnitsky, B.D. Rao, Sparse signal reconstruction from limited data using FOCUSS: A re-weighted minimum norm algorithm. *IEEE Trans. Signal Process.* **45**(3), 600–616 (1997)
12. M. Grant, S. Boyd, CVX: Matlab Software for Disciplined Convex Programming version 2.1. <http://cvxr.com/cvx>, Mar 2015
13. M. Grant, S. Boyd, Graph implementations for nonsmooth convex programs, in *Recent Advances in Learning and Control*, ed. by V. Blondel, S. Boyd, H. Kimura. Lecture Notes in Control and Information Sciences (Springer, Heidelberg, 2008), pp. 95–110
14. S. Ji, Y. Xue, L. Carin, Bayesian compressive sensing. *IEEE Trans. Signal Process.* **56**(6), 2346–2356 (2008)

15. S. Li, Concise formulas for the area and volume of a hyperspherical cap. *Asian J. Math. Stat.* **4**(1), 66–70 (2011)
16. D. Needell, J.A. Tropp, CoSaMP: Iterative signal recovery from incomplete and inaccurate samples. *Appl. Comput. Harmon. Anal.* **26**(3), 301–321 (2009)
17. S. Rangan, Generalized approximate message passing for estimation with random linear mixing, in *2011 IEEE International Symposium on Information Theory Proceedings*, IEEE, July 2011, pp. 2168–2172

Adapted Compressed Sensing for Effective Hardware  
Implementations

A Design Flow for Signal-Level Optimization of  
Compressed Sensing Stages

Mangia, M.; Pareschi, F.; Cambareri, V.; Rovatti, R.;  
Setti, G.

2018, XIV, 319 p. 180 illus., 142 illus. in color.,  
Hardcover

ISBN: 978-3-319-61372-7