

Sparse Reconstruction via Bayesian Model Averaging

Phil Schniter



July 2009

(Joint work with Dr. Lee Potter, Mr. Justin Ziniel, and Mr. Subhojit Som)

Sparse Reconstruction:

Estimate *sparse* x from the *under-determined* noisy linear mixture:

$$y = Ax + e \quad \text{for known } A \in \mathbb{C}^{M \times N}, \text{ with } M \ll N.$$

- *Under-determined* means that the number of measurements (M) is less than the number of unknowns (N).

In general, there is no unique solution!

- *Sparse* means that the number of nonzero coefs (K) is relatively few.

This may help to solve the problem. . .

But why do we care about this problem?

Standard Two-Step Data Acquisition:

1. Nyquist-rate sampling,
2. Lossy compression:
 - (a) take the discrete wavelet transform (DWT),
 - (b) keep only the large DWT coefficients.

Why? Because the DWT of a “structured” signal is sparse:

$$\boxed{\mathbf{W}z = \mathbf{x} + \mathbf{r}} : \begin{cases} \mathbf{W} & \text{unitary DWT operator,} \\ z & \text{structured signal,} \\ \mathbf{x} & \text{large wavelet coefficients. . . sparse,} \\ \mathbf{r} & \text{small residual. . . non-sparse.} \end{cases}$$

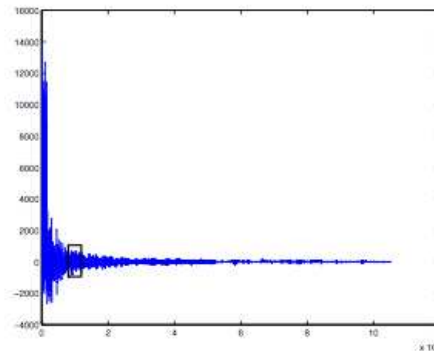
Furthermore, the DWT is *universal*; it doesn't need to know the particular “structure” of \mathbf{x} !

DWT Example:

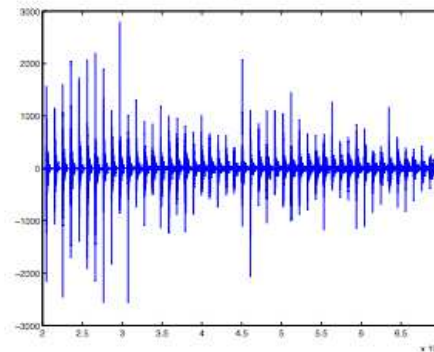
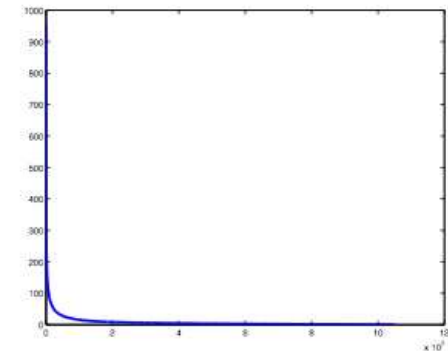


1 megapixel image

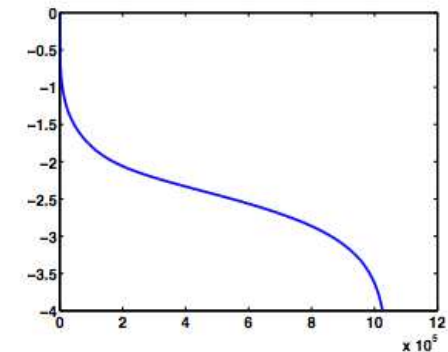
wavelet coeffs



(sorted)



zoom in



(\log_{10} sorted)

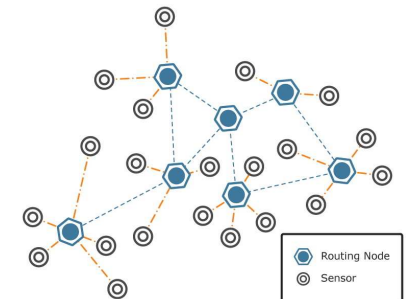
Typically: $\text{MSE} \approx -20$ dB from only 2.5% of DWT coefficients!

What's Wrong With Sampling-then-Compressing:

Sampling at the Nyquist rate produces many samples!

This poses problems when sampling is expensive, e.g.,

1. Magnetic resonance imaging (MRI): sampling is very time-intensive.
2. Sensor networks with “dumb” nodes: before compression, samples are communicated through a network.



Is there a more efficient way to sample?

An Alternative — Compressive Sampling/Sensing:

Main Idea:

For an N -dimensional signal z , take $M \ll N$ *linear* measurements

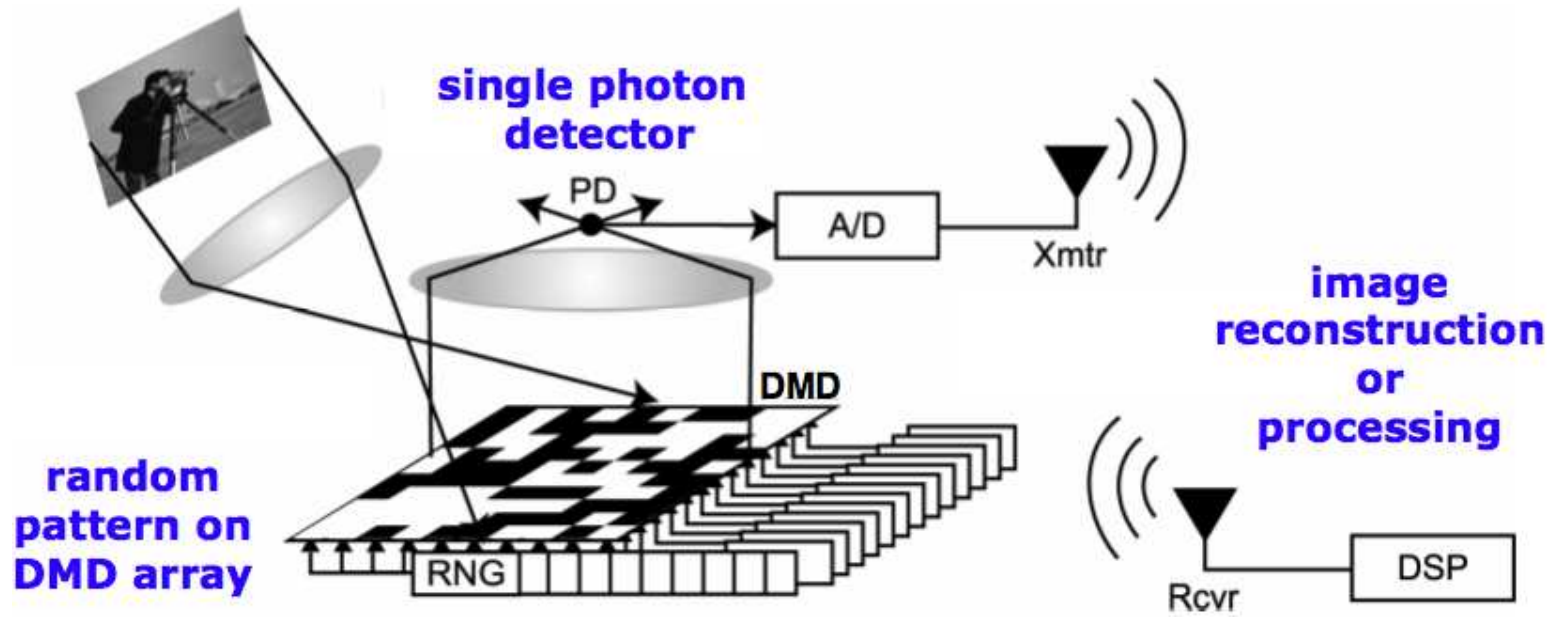
$y = \Phi z$ from which z can be accurately reconstructed.

How should we construct Φ ?

- Since we make no prior assumptions on the signal's eigen-structure, the measurement scheme Φ must be *universal*.
- Intuitively, measurements should collect approximately equal energy from *all* subspaces of \mathbb{C}^N .

Use a measurement matrix $\Phi \in \mathbb{C}^{M \times N}$ with random entries!

Example — Single-Pixel Camera (Rice University):



target
65536 pixels



11000 measurements
(16%)



1300 measurements
(2%)



This and the previous DWT figure courtesy of Rich Baraniuk, Rice University.

An Alternative — Compressive Sampling/Sensing:

How many measurements M do we need?

- Say z belongs to a class of signals with K dominant eigenspaces.
- These K subspaces can be arranged $\binom{N}{K}$ ways within the N dimensional space, so the signal z carries at least

$$\log_2 \binom{N}{K} \geq K \log_2 \left(\frac{N}{K} \right) \text{ bits of information.}$$

- For measurements $\mathbf{y} = \Phi \mathbf{z} + \mathbf{v}$ in AWGN, we learn at most

$$\frac{1}{2} \log_2(1 + \text{SNR}) \text{ bits per measurement,}$$

suggesting that we need at least

$$M \geq \frac{2}{\log_2(1 + \text{SNR})} K \log_2 \left(\frac{N}{K} \right) \text{ measurements.}$$

Compressive Sensing:

Essential components:

Measurement: $\mathbf{y} = \Phi \mathbf{z} + \boldsymbol{\nu}$ for $\Phi \in \mathbb{C}^{M \times N}$ and noise $\boldsymbol{\nu}$.

Compressibility: $\mathbf{W} \mathbf{z} = \mathbf{x} + \mathbf{r}$ for K -sparse \mathbf{x} , small residual \mathbf{r} , and unitary transform \mathbf{W} .

Putting these together, we get

$$\begin{aligned} \mathbf{y} &= \underbrace{\Phi \mathbf{W}^H}_{\mathbf{A}} \mathbf{x} + \underbrace{\Phi \mathbf{W}^H \mathbf{r} + \boldsymbol{\nu}}_{\mathbf{e}} \\ &= \mathbf{A} \mathbf{x} + \mathbf{e} \end{aligned}$$

where $\mathbf{A} \in \mathbb{C}^{M \times N}$ and $K < M \ll N$.

The remaining challenge:

Reconstruct the sparse signal representation \mathbf{x} from the “compressed” measurements $\mathbf{y} = \mathbf{A} \mathbf{x} + \mathbf{e}$.

Other Applications of Sparse Reconstruction:

Sparse channel estimation:

$$\mathbf{r} = \mathbf{A}\mathbf{x} + \mathbf{e} \quad \text{for} \quad \begin{cases} \mathbf{x}: \text{ sparse channel impulse response,} \\ \mathbf{A}: \text{ pilot symbol matrix.} \end{cases}$$

- When the problem is under-determined ($M < N$):
 - \rightsquigarrow *Sparsity is **needed** to solve the problem!*
 - Span of \mathbf{r} limited to ensure that channel is time-invariant over block.
 - Span of \mathbf{r} limited due to small number of pilot symbols.

- When the problem is not under-determined ($M \geq N$):

\rightsquigarrow *Sparsity can be leveraged to **improve** performance!*

$$\text{since } \mathbb{E}\{\|\hat{\mathbf{x}}_{\text{LS}} - \mathbf{x}\|_2^2\} \approx \frac{K\sigma_e^2}{M\sigma_a^2} \quad \text{for } K = \underbrace{\|\mathbf{x}\|_0}_{\# \text{ nonzero coefs}}.$$

Solving the Sparse Reconstruction Problem:

Key question:

How do we use the sparsity of \mathbf{x} to solve the noisy under-determined inverse problem $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{e}$?

Popular Techniques:

1. Convex optimization of constrained ℓ_1 criteria.
2. Matching pursuits.
3. Bayesian approaches.

Notation: $\|\mathbf{x}\|_p = \sqrt[p]{\sum_n |x_n|^p}$ is the " ℓ_p norm."

The Canonical “Sparse-Approximation” Problem:

Find the sparsest x which explains y up to a specified tolerance of ϵ :

$$\hat{x} = \arg \min_x \underbrace{\|x\|_0}_{\# \text{ nonzero coefs}} \text{ s.t. } \|y - Ax\|_2 \leq \epsilon.$$

Key points:

1. NP-hard: need to check all configurations of non-zero coefficients!
2. May not be the “right” objective.

Let’s think about this sparse-approximation problem geometrically...

A Toy Example:

Consider the setup

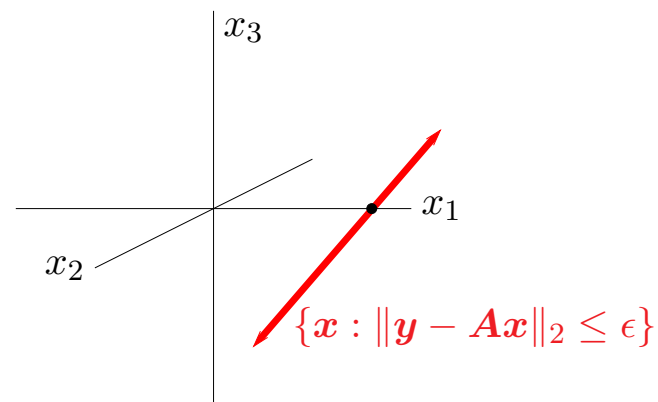
$$\begin{bmatrix} \bullet \\ \bullet \end{bmatrix} = \begin{bmatrix} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{bmatrix} \begin{bmatrix} \bullet \\ \bullet \\ \bullet \end{bmatrix} + \begin{bmatrix} \bullet \\ \bullet \end{bmatrix} \quad \begin{array}{l} N = 3 \\ M = 2 \\ K = 1 \end{array}$$

Since $N = M + 1$,

- the set $\{x : y = Ax\}$ is described by a line (via $\text{null}(A)$), and
- the set $\{x : \|y - Ax\|_2 \leq \epsilon\}$ is described by an ϵ -thin rod.

Since $K = 1$,

- the true x intersects one of the coordinate axes. (But which one?)

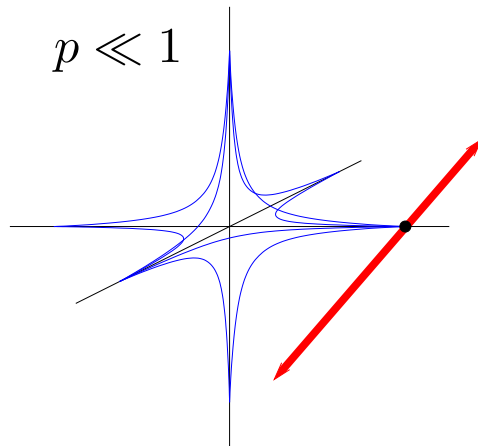


The Geometry of Constrained ℓ_p -Minimization:

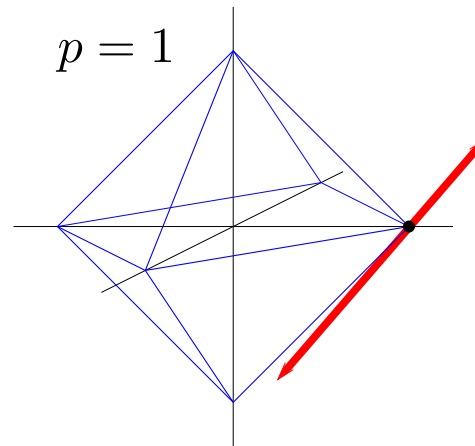
Now consider, for some general $p > 0$, the optimization problem:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{x}\|_p \quad \text{s.t.} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2 \leq \epsilon.$$

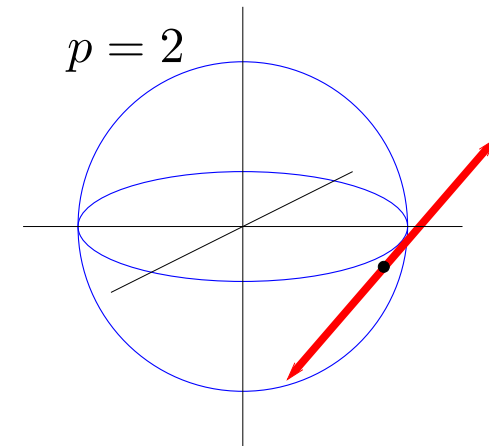
$\hat{\mathbf{x}}$ can be found by growing the ℓ_p -ball until it touches the ϵ -rod:



Solution definitely sparse
but problem is **NP hard**.



Solution usually sparse
and problem is **convex**!



Solution is **not sparse**;
 \Leftrightarrow LS when $\epsilon = 0$.

This suggests to use the ℓ_1 norm as a surrogate for the ℓ_0 norm.

1) Constrained ℓ_1 -Minimization:

For the constrained- ℓ_1 approach (known as “LASSO”)

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{x}\|_1 \text{ s.t. } \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2 \leq \epsilon,$$

there exist elegant theorems which say that, given

- enough measurements (e.g., $M \gtrsim K \log(N - K)$) and
- sufficiently well-behaved \mathbf{A} (e.g., nearly uncorrelated columns),

$\|\hat{\mathbf{x}} - \mathbf{x}\|_2$ will be very small with very high probability.

[1] D. L. Donoho, M. Elad, and V. N. Temlyakov, “Stable recovery of sparse overcomplete representations in the presence of noise,” *IEEE Trans. Inform. Theory*, vol. 52, no. 1, 2006.

[2] J. A. Tropp, “Just relax: Convex programming methods for identifying sparse signal,” *IEEE Trans. Info. Theory*, vol. 51, no. 3, 2006.

[3] E. Candès, J. Romberg, and T. Tao, “Stable signal recovery from incomplete and inaccurate measurements,” *Commun. on Pure and Applied Math.*, vol. 59, no. 8, 2006.

*But, \mathbf{A} may not be well-behaved, especially when M and N are not huge!
Also, generally incompatible with complex-valued x .*

2) Matching-Pursuit Algorithms:

- Basic “matching pursuit” algorithm:

Similar to “successive interference cancellation” for CDMA:

1. Find the column \mathbf{a}_i of \mathbf{A} that is most correlated with \mathbf{y} .
 2. Estimate the corresponding signal coefficient x_i using least-squares.
 3. Compute the residual: $\mathbf{r} = \mathbf{y} - \mathbf{a}_i \hat{x}_i$.
 4. Repeat with \mathbf{r} in place of \mathbf{y} .
- More sophisticated versions, like “orthogonal matching pursuit” (OMP), are more robust to correlation among the columns of \mathbf{A} .

Matching-Pursuit Theory:

For some matching-pursuit algorithms, one can prove that, with

- enough measurements (i.e., $M \gtrsim K \log N$), and
- sufficiently well-behaved \mathbf{A} (e.g., nearly orthogonal columns),

reconstruction error is small with high probability.

[1] J. A. Tropp and A. C. Gilbert, "Signal Recovery From Random Measurements Via Orthogonal Matching Pursuit," *IEEE Trans. Inform. Thy.*, Dec. 2007.

[2] D. Needell and J. A. Tropp, "CoSaMP: Iterative signal recovery from incomplete and inaccurate samples," *Harmon. Anal.*, 2008.

[3] W. Dai and O. Milenkovic, "Subspace pursuit for compressive sensing: Closing the gap between performance and complexity," Preprint 2008.

Practical, but how well does it work when \mathbf{A} is not well-behaved?

3) Bayesian Approaches to Sparse Reconstruction:

Say that we have some prior statistical knowledge of

- the pattern of active coefficients,
- the values of active coefficients,
- the noise.

Can we take advantage of this knowledge for sparse reconstruction?

Three of the most popular Bayesian strategies are...

- a) Laplacian signal prior,
- b) The relevance vector machine.
- c) **Bayesian variable selection and Bayesian model averaging,**

3a) Laplacian Signal Prior:

- If we assume σ^2 -variance AWGN and signal \mathbf{x} such that

$$p(\mathbf{x}) \propto e^{-\tau \|\mathbf{x}\|_p^p},$$

then the MAP estimate becomes

$$\begin{aligned} \hat{\mathbf{x}} &= \arg \max_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}) \\ &= \arg \min_{\mathbf{x}} \sigma^{-2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \tau \|\mathbf{x}\|_p^p \\ &= \arg \min_{\mathbf{x}} \|\mathbf{x}\|_p \quad \text{s.t.} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2 \leq \epsilon(\tau\sigma^2) \end{aligned}$$

which is the constrained ℓ_p -optimization problem we saw earlier.

- Choosing the Laplacian prior (i.e., $p = 1$), we know that \mathbf{x} will be sparse and that $\hat{\mathbf{x}}$ can be obtained via convex programming.

Interesting, but... Physical meaning of Laplace prior? Choice of τ ?

3b) The Relevance Vector Machine (RVM):

- To model coefficient activity, use “precisions” $\alpha \in (\mathbb{R}^+)^N$:

$$\mathbf{x}|\alpha \sim \prod_n \mathcal{N}(0, \alpha_n^{-1}) \quad \text{and} \quad \alpha \sim \text{iid } \Gamma(0, 0)$$

$$\mathbf{e}|\beta \sim \prod_n \mathcal{N}(0, \beta^{-1}) \quad \text{and} \quad \beta \sim \Gamma(0, 0)$$

As $\alpha_n \rightarrow \infty$, the coefficient x_n is effectively “turned off”.

- The use of gamma hyperpriors leads to the convenient posterior

$$p(\mathbf{x}|\mathbf{y}, \alpha, \beta) \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad \text{for} \quad \begin{cases} \boldsymbol{\mu} = \beta \boldsymbol{\Sigma} \mathbf{A}^T \mathbf{y} \\ \boldsymbol{\Sigma} = (\beta \mathbf{A}^T \mathbf{A} + \mathcal{D}(\alpha))^{-1} \end{cases}$$

and thus the convenient estimate $\hat{\mathbf{x}}_{\text{MMSE}} = \boldsymbol{\mu}$.

- The EM algorithm can be used to estimate $\{\alpha, \beta\}$ jointly with $\{\boldsymbol{\mu}, \boldsymbol{\Sigma}\}$.

Can implement with an $\mathcal{O}(NK^2)$ recursion after an $\mathcal{O}(N^2M)$ initialization.

[1] Tipping, “Sparse Bayesian learning and the relevance vector machine,” *J. Machine Learning Res.*, 2001.

[2] Wipf and Rao, “Sparse Bayesian learning for basis selection,” *IEEE Trans. Signal Processing*, 2004.

[3] Ji, Xue, and Carin, “Bayesian compressive sensing,” *IEEE Trans. Signal Processing*, 2008.

3c) Variable Selection:

- Using S to denote the set of active-coefficient indices, we can write

$$\mathbf{y} = \mathbf{A}_S \mathbf{x}_S + \mathbf{e}.$$

- With S known, estimation of the nonzero coefficients \mathbf{x}_S is easy:

$$\begin{aligned}\hat{\mathbf{x}}_{\text{LS}|S} &= (\mathbf{A}_S^T \mathbf{A}_S)^{-1} \mathbf{A}_S^T \mathbf{y} \\ \hat{\mathbf{x}}_{\text{MMSE}|S} &= (\mathbf{A}_S^T \mathbf{A}_S + \sigma_e^2 \mathbf{I})^{-1} \mathbf{A}_S^T \mathbf{y}\end{aligned}$$

This motivates the variable-selection problem:

From $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{e}$, estimate the active-coefficient set S .

a long-standing problem in statistics!

[1] Hocking, “The analysis and selection of variables in linear regression,” *Biometrics*, 1976.

Non-Bayesian Variable Selection:

Consider the GLRT detector:

$$\begin{aligned}
 \hat{S}_{\text{GLRT}} &= \arg \max_S p(\mathbf{y}|S, \hat{\mathbf{x}}_{\text{ML}|S}) \\
 &= \arg \min_S \|\mathbf{y} - \mathbf{A}_S \underbrace{(\mathbf{A}_S^T \mathbf{A}_S)^{-1} \mathbf{A}_S^T \mathbf{y}}_{\hat{\mathbf{x}}_{\text{ML}|S}}\|_2^2 \\
 &= \arg \min_S \|\mathcal{P}_{\mathbf{A}_S}^\perp \mathbf{y}\|_2^2 \\
 &= \text{any } S : |S| \geq M.
 \end{aligned}$$

It fails!

This happens whenever the models are “nested”.

Bayesian Variable Selection:

Consider the MAP model estimate:

$$\begin{aligned}\hat{S}_{\text{MAP}} &= \arg \max_S p(S|\mathbf{y}) \\ &= \arg \max_S p(\mathbf{y}|S)p(S) \\ &= \arg \max_S \int_{\mathcal{X}} \underbrace{p(\mathbf{y}|S, \mathbf{x})}_{\mathcal{N}} p(\mathbf{x}|S) d\mathbf{x} \cdot p(S).\end{aligned}$$

Need to specify the priors $p(\mathbf{x}|S)$ and $p(S)$...

- [1] Lempers, *Posterior probabilities of alternative linear models*, Rotterdam: Rotterdam Univ. Press, 1971
- [2] Mitchell & Beauchamp, "Bayesian variable selection in linear regression," *J. Amer. Statist. Assoc.*, 1988.
- [3] George & McCulloch, "Variable selection via Gibbs sampling," *J. Amer. Statist. Assoc.*, 1993.
- [4] Smith & Kohn, "Nonparametric regression using Bayesian variable selection," *J. Econometrics*, 1996.
- [5] George & McCulloch, "Approaches for Bayesian variable selection," *Statist. Sinica*, 1997.
- [6] George, "The variable selection problem," *J. Amer. Statist. Assoc.*, 2000.

Popular BVS Priors:

- iid Bernoulli coefficient-activity:

$$p(S) = \lambda^{|S|}(1 - \lambda)^{(N-|S|)} \quad \text{where } \lambda < 0.5 \text{ induces sparsity,}$$

- Gaussian active-coefficients \mathbf{x}_S :

$$p(\mathbf{x}_S|S) \sim \mathcal{N}(\mu\mathbf{1}_{|S|}, \mathbf{R}_S)$$

$$\text{for } \begin{cases} \mathbf{R}_S = \sigma_x^2 \mathbf{I}_{|S|}, & \mu \in \mathbb{R} & \text{“iid”} \\ \mathbf{R}_S = \sigma_x^2 (\mathbf{A}_S^T \mathbf{A}_S)^{-1}, & \mu = 0 & \text{“Zellner } g\text{-prior”} \end{cases}$$

where the hyperparameters $\{\lambda, \mu, \sigma_x^2, \sigma_e^2\}$ could be treated as...

1. *random*: assign non-informative conjugate priors & integrate out unknowns,
2. *deterministic*: use the EM-algorithm to estimate hyperparameters.

[1] Cui & George, “Empirical Bayes vs. fully Bayes variable selection,” *J. Statist. Planning Infer.*, 2008.

BVS Posteriors:

Fixing $\{\lambda, \mu, \sigma_x^2, \sigma_e^2\}$, we get

- the model's posterior log-density:

$$\ln p(S|\mathbf{y}) = -\frac{1}{2} \|\mathbf{y} - \mu \mathbf{A}_S \mathbf{1}_{|S|}\|_{\Phi_S^{-1}}^2 - \frac{1}{2} \ln \det(\Phi_S) - |S| \ln\left(\frac{1-\lambda}{\lambda}\right) + C,$$

where Φ_S denotes the observation covariance matrix given S :

$$\Phi_S = \begin{cases} \sigma_x^2 \mathbf{A}_S \mathbf{A}_S^T + \sigma_e^2 \mathbf{I}_{|S|} & \text{(iid)} \\ \sigma_x^2 \mathbf{A}_S (\mathbf{A}_S^T \mathbf{A}_S)^{-1} \mathbf{A}_S^T + \sigma_e^2 \mathbf{I}_{|S|} & \text{(Zellner)} \end{cases}.$$

- the (S -conditional) coefficient density:

$$p(\mathbf{x}_S | \mathbf{y}, S) \sim \mathcal{N}(\boldsymbol{\mu}_S, \boldsymbol{\Sigma}_S)$$

where

$$\begin{aligned} \boldsymbol{\mu}_S &= \mu \mathbf{1}_{|S|} + \mathbf{R}_S \mathbf{A}_S^T \Phi_S^{-1} (\mathbf{y} - \mu \mathbf{A}_S \mathbf{1}_{|S|}) = \hat{\mathbf{x}}_{\text{MMSE}|S} \\ \boldsymbol{\Sigma}_S &= \mathbf{R}_S - \mathbf{R}_S \mathbf{A}_S^T \Phi_S^{-1} \mathbf{A}_S \mathbf{R}_S = \text{cov}(\hat{\mathbf{x}}_{\text{MMSE}|S}). \end{aligned}$$

Connections to Model-Order Selection:

Under the Zellner prior, it can be shown that

$$\hat{S}_{\text{MAP}} = \arg \min_S \left\{ \frac{1}{\sigma_e^2} \|\mathbf{y} - \mathbf{A}_S \hat{\mathbf{x}}_{\text{LS}|S}\|_2^2 + |S| \cdot \eta \right\}$$

$$\text{for } \eta = \frac{\sigma_x^2 + \sigma_e^2}{\sigma_x^2} \ln \left(\left(1 + \frac{\sigma_x^2}{\sigma_e^2}\right) \left(\frac{1-\lambda}{\lambda}\right)^2 \right).$$

Note the close connections to “information theoretic” model-order selectors:

$$\hat{S}_{\text{AIC}} = \arg \min_S \left\{ \frac{1}{\sigma_e^2} \|\mathbf{y} - \mathbf{A}_S \hat{\mathbf{x}}_{\text{LS}|S}\|_2^2 + |S| \cdot 2 \right\}$$

$$\hat{S}_{\text{BIC}} = \arg \min_S \left\{ \frac{1}{\sigma_e^2} \|\mathbf{y} - \mathbf{A}_S \hat{\mathbf{x}}_{\text{LS}|S}\|_2^2 + |S| \cdot \ln M \right\}$$

$$\hat{S}_{\text{RIC}} = \arg \min_S \left\{ \frac{1}{\sigma_e^2} \|\mathbf{y} - \mathbf{A}_S \hat{\mathbf{x}}_{\text{LS}|S}\|_2^2 + |S| \cdot 2 \ln N \right\}.$$

[1] George & Foster, “Calibration and empirical Bayes variable selection,” *Biometrika*, 2000.

Connections to Noncoherent Decoding:

- Consider a generic communication system with vectorized observations

$$\mathbf{y} = \mathbf{A}_j \mathbf{h} + \mathbf{e},$$

where $j =$ codeword index, $\mathbf{h} =$ channel gains, and $\mathbf{e} =$ AWGN.

- In *non-coherent decoding*, \mathbf{h} is known only in distribution, e.g., $\mathbf{h} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}_N)$ for Rayleigh fading. Then

$$\hat{j}_{\text{MAP}} = \arg \max_{j=1, \dots, J} p(j|\mathbf{y}).$$

- This can be interpreted in our sparse reconstruction framework via

$$\mathbf{A} \triangleq [\mathbf{A}_1 \cdots \mathbf{A}_J] \quad \text{and} \quad \mathbf{x} \triangleq [\mathbf{0}^T, \dots, \mathbf{h}^T, \dots, \mathbf{0}^T]^T$$

with a constraint on the admissible models $S \in \mathcal{S}$:

$$\mathcal{S} \triangleq \{(1, \dots, K), (K+1, \dots, 2K), \dots, (JK-K+1, \dots, JK)\},$$

with the result that

$$\hat{S}_{\text{MAP}} = \arg \max_{S \in \mathcal{S}} p(S|\mathbf{y}) \Leftrightarrow \hat{j}_{\text{MAP}}$$

Bayesian Model Averaging:

- Previously we motivated Bayesian variable selection, e.g.,

$$\hat{S}_{\text{MAP}} = \arg \max_S p(S|\mathbf{y})$$

for subsequent use in a *conditional* estimation strategy, e.g.,

$$\hat{\mathbf{x}}_{\text{MMSE}|\hat{S}_{\text{MAP}}} = \mathbb{E}\{\mathbf{x}|\mathbf{y}, \hat{S}_{\text{MAP}}\}.$$

- But having access to the “soft information” $\{p(S|\mathbf{y})\}_{\forall S}$ allows more sophisticated *unconditional* estimates, e.g.,

$$\hat{\mathbf{x}}_{\text{MMSE}} = \sum_S \hat{\mathbf{x}}_{\text{MMSE}|S} p(S|\mathbf{y})$$

that are well approximated by summing over the *few* most probable S .

[1] Leamer, *Specification Searches*, New York: Wiley 1978.

[2] Raftery, Madigan, & Hoeting, “Bayesian model averaging for linear regression models,” *J. Amer. Statist. Assoc.*, 1997.

[3] Clyde and George, “Model Uncertainty,” *Statist. Sci.*, 2004.

Implementation of BVS/BMA:

- The conventional approach to determining the set of high-probability S is based on random search (e.g., Gibbs Sampling or Markov Chain Monte Carlo).
- Recently, a (non-exhaustive) tree search has been proposed to learn the set of high-probability S (and their $\hat{\mathbf{x}}_{\text{MMSE}|S}$) with very low complexity:
 - iid Gaussian \mathbf{x}_S : “Fast Bayesian matching pursuit”
 - Zellner Gaussian \mathbf{x}_S : “Fast Zellner matching pursuit”

In fact, the complexity order equals that of OMP: $\mathcal{O}(MNK)$.

- An expectation/maximization (EM) approach can be used to optimally tune the hyperparameters $\{\lambda, \mu, \sigma_x^2, \sigma_e^2\}$ wrt the data.

[1] Schniter, Potter, and Ziniel, “Fast Bayesian matching pursuit,” *ITA*, 2008.

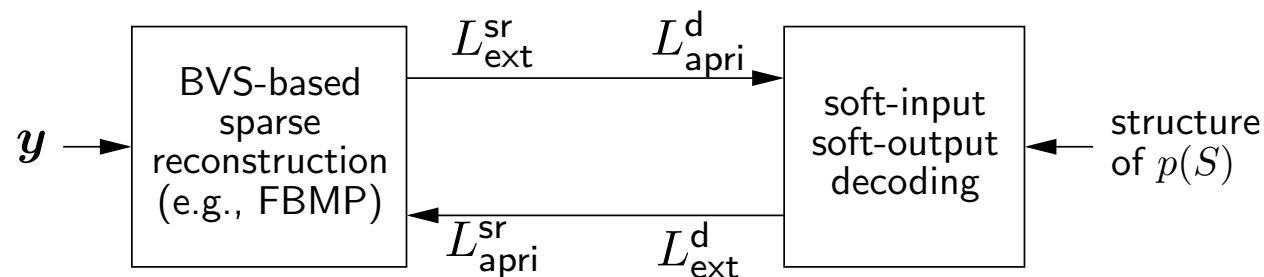
[2] Schniter, Potter, and Ziniel, “Fast Bayesian matching pursuit: Model uncertainty and parameter estimation for sparse linear models,” Preprint, 2008.

Turbo Implementation:

Recall that

- BVS is based on *binary* indicators of coefficient activity.
- BVS is a *soft-input soft-output* form of sparse reconstruction (i.e., it takes $p(S)$ as input and generates $p(S|\mathbf{y})$ as output).

With complicated $p(S)$, exact BVS implication may be difficult. Thus one might consider an iterative approach:



that passes extrinsic log-likelihood ratios on the indicator variables.

[1] Schniter, "Bayesian Sparse Reconstruction and Tracking using the Turbo Principle," *Preprint*, 2009.

BMA versus RVM:

- To parameterize coefficient activity, RVM uses the *continuous* variables α , while BMA uses the *discrete* set S .
- RVM has trade-off parameters that must be tuned using *cross-validation*, whereas the BMA hyperparameters can be tuned via the *EM algorithm*.
- RVM infers *marginal* coefficient activity, whereas BMA infers *joint* coefficient activity.
- Upon termination, the RVM posterior is Gaussian

$$p(\mathbf{x}|\mathbf{y}) \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

whereas the BMA posterior is a Gaussian mixture:

$$p(\mathbf{x}|\mathbf{y}) \sim \sum_S \mathcal{N}(\boldsymbol{\mu}_S, \boldsymbol{\Sigma}_S) p(S|\mathbf{y}) \quad \dots \text{more informative}$$

- Fast implementations of RVM and BMA have roughly the same complexity.
- Simulation results suggest better performance for BMA.

Numerical Experiments — “Compressible” Signal:

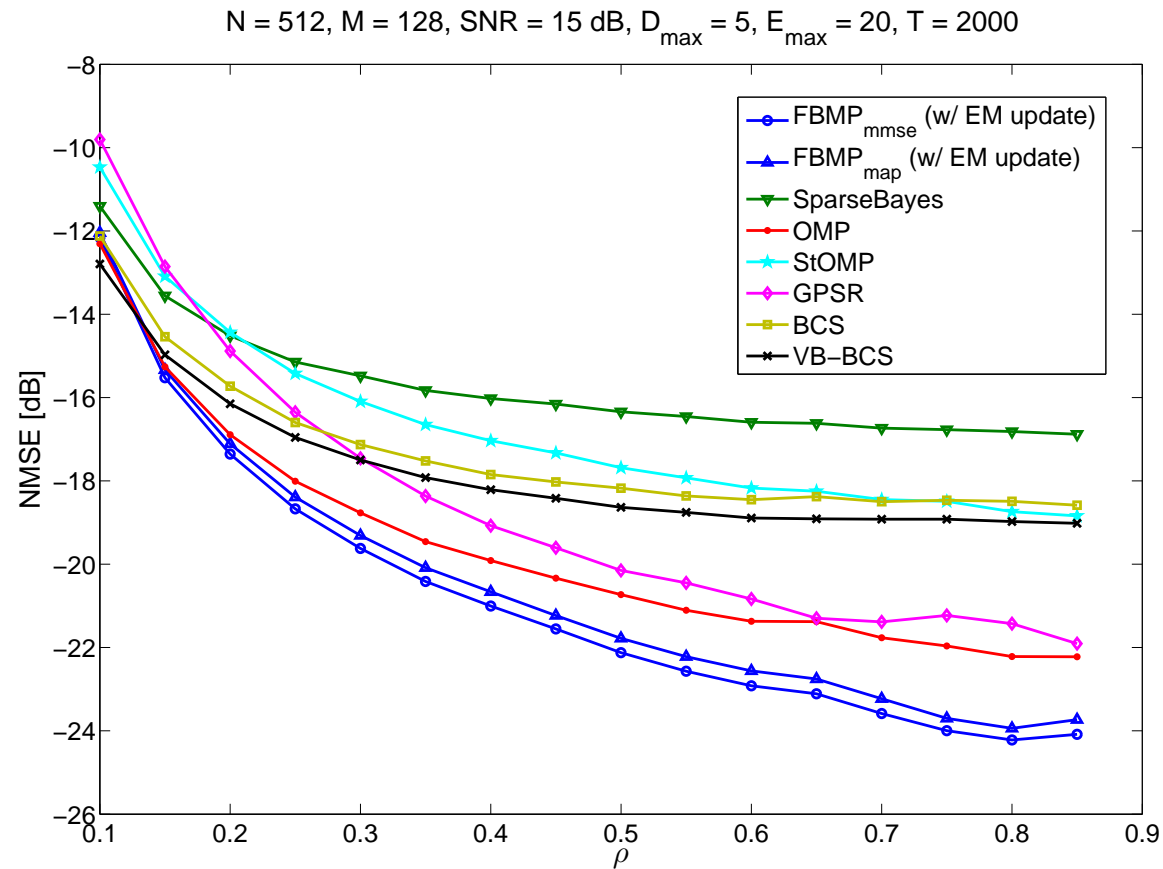
Setup: $N = 512$
 $M = 128$
 \mathbf{A} : i.i.d. $\mathcal{N}(0, 1)$ with columns scaled to unit norm
 \mathbf{x} : $x_n = e^{-\rho n}$ (flipped/shuffled) for decay rate $\rho \in (0, 1)$
 SNR = 15dB

Algorithms:

OMP – Tropp & Gilbert
 StOMP – Donoho, Tsaig, Drori & Starck
 GPSR-Basic – Figueiredo, Nowak & Wright ($\min_{\mathbf{x}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \tau\|\mathbf{x}\|_1$)
 SparseBayes – Wipf & Rao (RVM)
 BCS – Ji & Carin (RVM)
 FBMP – Schniter, Potter & Ziniel (BMA)

Performance: $\text{NMSE} \triangleq \text{Avg} \left\{ \frac{\|\hat{\mathbf{x}} - \mathbf{x}\|_2^2}{\|\mathbf{x}\|_2^2} \right\}$ over 2500 random trials.

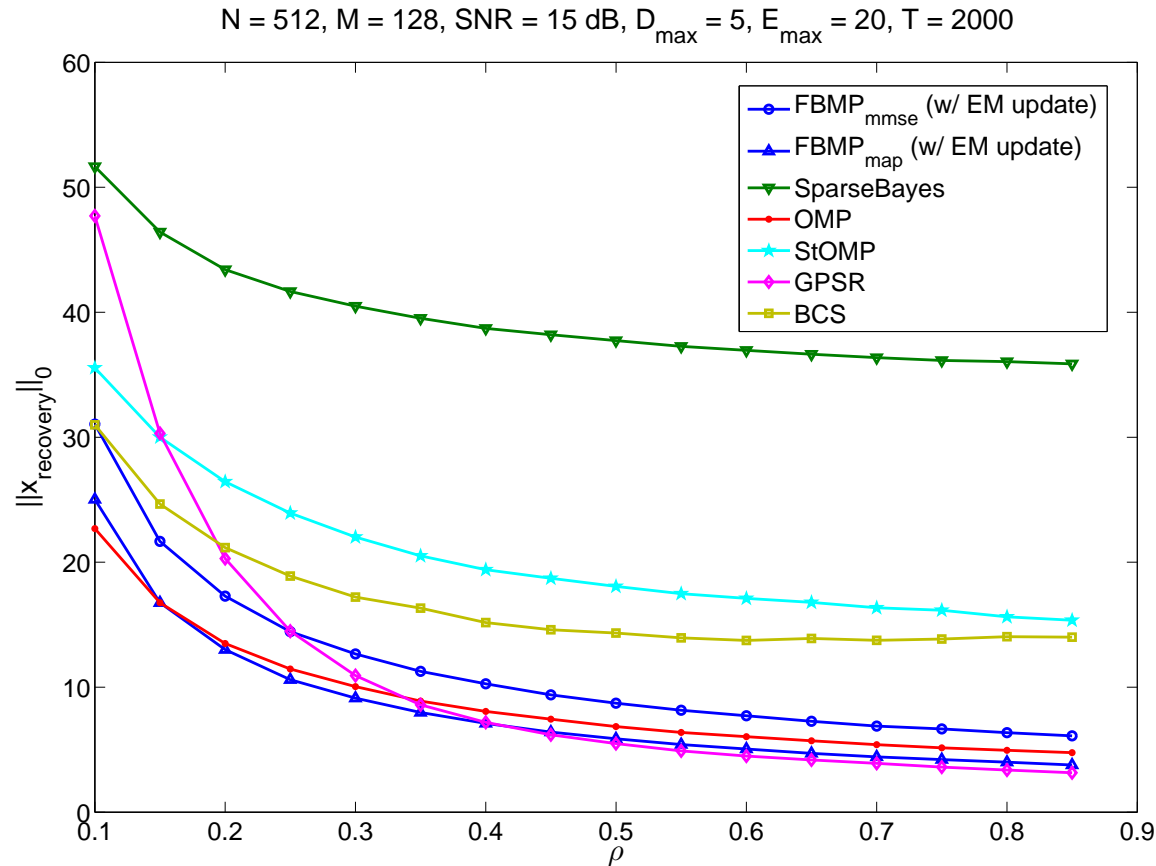
NMSE versus decay rate ρ :



FBMP outperformed GPSR and OMP by 2 dB and others by much more.

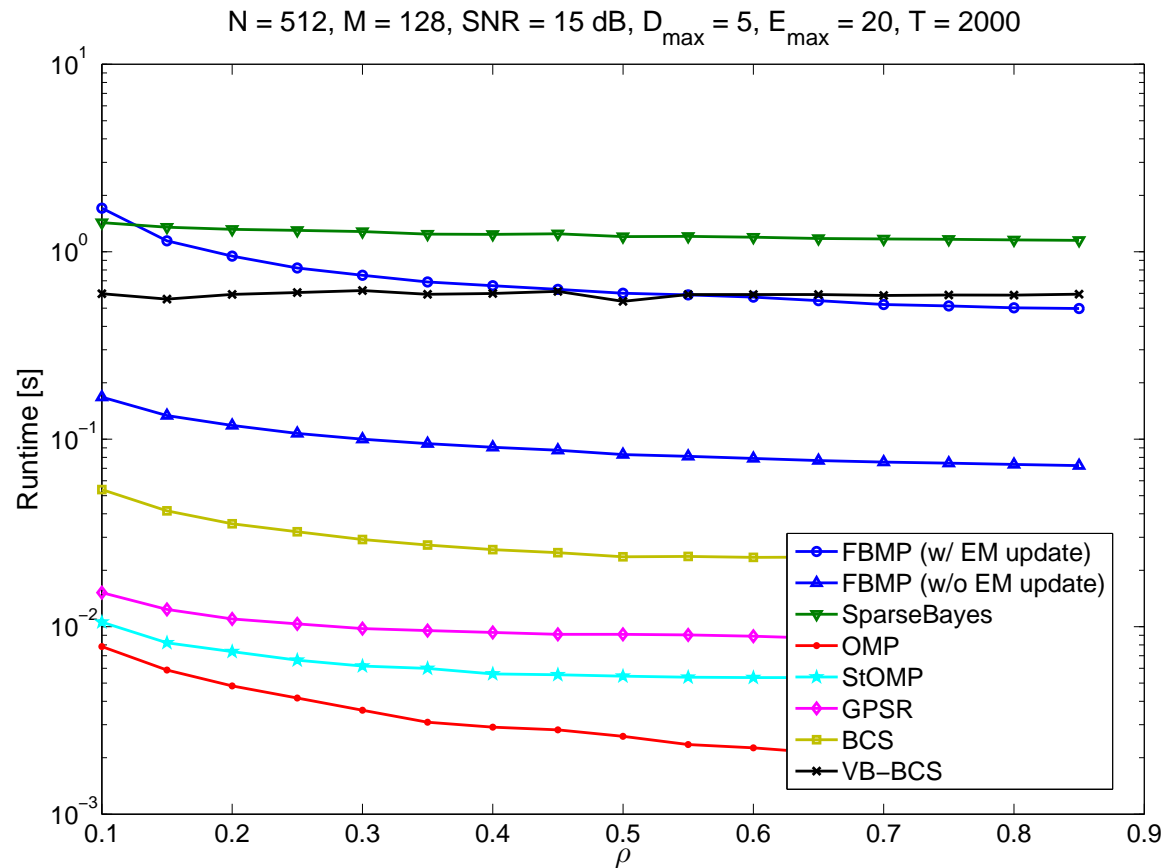
Note: The signal priors favor GPSR!

Sparsity of estimate versus decay rate ρ :



The estimates returned by FBMP are among the sparsest.
 (While BMA is generally not sparse, FBMP is, due to non-exhaustive search.)

Runtime versus decay rate ρ :



FBMP is on par with other Bayesian algorithms, but slower than OMP and convex programming algorithms.

Performance Guarantees for MAP Variable Selection:

Under the iid-Beroulli/Gaussian signal model and a Restricted Isometry Property (RIP) for \mathbf{A} :

There exists $\delta_K \in (0, 1)$ such that

$$\forall \mathbf{x} \text{ s.t. } \|\mathbf{x}\|_0 = K, \quad (1 - \delta_K)\|\mathbf{x}\|_2^2 \leq \|\mathbf{A}\mathbf{x}\|_2^2 \leq (1 + \delta_K)\|\mathbf{x}\|_2^2,$$

it has been recently shown that the following properties hold *with high probability* for reasonably small constants K_1, K_2, K_3, K_4 :

1. The energy of the missed signal coefficients is upper bounded by $K_1 M \sigma_e^2$.
2. No active coefficients are missed when $|\mu| > 4\sigma_1 + K_2 \sqrt{M} \sigma_e^2$.
3. No coefficients are falsely detected when $|\mu| > K_3 \sqrt{M} \sigma_1 + K_4 \sqrt{M} \sigma_e^2$.

[1] Som and Potter, "MAP Estimation of Support Set in Sparse Signal Recovery," *Preprint*, 2009.

PWEP of Model Selection under general “Information Criteria”:

Lemma 1 For generic \mathcal{S} , the pairwise error probability of

$$\hat{S} = \arg \min_{S \in \mathcal{S}} \left\{ \frac{1}{\sigma_e^2} \|\mathbf{y} - \mathbf{A}_S \hat{\mathbf{x}}_{\text{LS}|S}\|_2^2 + \eta |S| \right\} \quad \text{under } \mathbf{x}_S | S \sim \mathcal{N}(\mathbf{0}, \gamma \sigma_e^2 \mathbf{I}_{|S|})$$

has the upper bound (tight as $\gamma \rightarrow \infty$):

$$P_{\hat{S}|S} \leq (\alpha_{\hat{S},S} \gamma)^{-K_m} C_{K_m, K_f}(\eta),$$

where K_m and K_f denote the # of missed and false-alarm coefficients, and

$$C_{K_m, K_f}(\eta) = \begin{cases} e^{(K_m - K_f)\eta} \sum_{k=0}^{K_f-1} \frac{(K_f - K_m)^k \eta^k}{k!} \binom{K_m + K_f - 1 - k}{K_m} & K_m \leq K_f, \\ \sum_{k=0}^{K_m} \frac{(K_m - K_f)^k \eta^k}{k!} \binom{K_m + K_f - 1 - k}{K_f - 1} & K_m > K_f. \end{cases}$$

$$\alpha_{\hat{S},S} = \lambda_{\min}(\mathbf{A}_m^T \mathbf{\Pi}_{\mathbf{A}_{\hat{S}}}^{\perp} \mathbf{A}_m) \quad \dots \text{Restricted Isometry Property}$$

[1] Schniter, “On Model Selection under the Generalized Information Criteria,” *Preprint*, 2009.

Conclusions:

- Bayesian variable selection (BVS) and Bayesian model averaging (BMA) are well-established statistical methods for sparse reconstruction.
- While BVS & BMA were previously considered to be “too expensive,” modern tree-search implementations have reasonable complexity.
- There are close connections between BVS and AIC/BIC/RIC, as well as between BVS and noncoherent decoding.
- Numerical experiments suggest that BMA yields NMSE superior to that of other state-of-the-art algorithms.
- Current work includes BVS/BMA performance analysis, turbo implementation, and applications in
 - medical imaging,
 - through-wall radar,
 - underwater acoustic channel tracking,
 - decoding with intermittent and degraded side-information.