

Fast Bayesian Matching Pursuit

Phil Schniter, Lee Potter, and Justin Ziniel



January 2008

Introduction:

The linear regression problem

“estimate sparse \mathbf{x} from measurements $\mathbf{y} = \mathbf{A}\mathbf{x} + \boldsymbol{\nu}$ ”

is often posed as a penalized least-squares (LS) problem:

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmin}} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{x}\|_p^p$$

$$\text{or } \hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmin}} \|\mathbf{x}\|_p^p \text{ s.t. } \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 < \epsilon$$

since the choice $p = 1$ provides a unique solution to the non-convex task

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmin}} \|\mathbf{x}\|_0 \text{ s.t. } \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 < \epsilon.$$

The LS approach can be interpreted as seeking the Bayesian MAP estimate of \mathbf{x} under the sparsity-inducing prior

$$p(\mathbf{x}) \sim \exp \left\{ -\frac{\lambda}{2} \|\mathbf{x}\|_p^p \right\}.$$

“Sparse Bayesian Learning” :

The method of “sparse Bayesian learning” explicitly adopts a Bayesian framework in which $\{x_i\}$ are independent such that

$$x_i \sim \mathcal{N}(0, \sigma_i^2)$$

σ_i^2 : unknown variance with a Gamma conjugate prior.

The EM algorithm is then used to find the MAP estimate.

However, the physical interpretation of the prior is not clear...

For example,

M. E. Tipping “Sparse Bayesian learning and relevance vector machine,” *J. Machine Learning Res.*, 2001.

D. Wipf and B. Rao, “Sparse Bayesian learning for basis selection,” *IEEE TSP*, 2004.

Our Problem Formulation:

Measurements:

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \boldsymbol{\nu} \in \mathbb{R}^M$$

\mathbf{A} : known mixing matrix

$\boldsymbol{\nu}$: AWGN with variance σ^2

\mathbf{x} : unknown sparse signal $\in \mathbb{R}^N$, (typically $N \gg M$).

Sparse Signal:

$$x_n | s_n \sim \begin{cases} \mathcal{N}(0, 1) & s_n = 1 \\ 0 & s_n = 0 \end{cases}$$

$s_n \sim \text{Bernoulli}(p_1)$ where typically $p_1 \ll 1$

$\{x_n\}_{n=0}^{N-1}, \{s_n\}_{n=0}^{N-1}$: i.i.d.

Reminiscent of the model adopted in

E. Larsson and Y. Selén, "Linear regression with a sparse parameter vector,"
IEEE TSP, Feb. 2007.

Objectives:

1. Basis Estimation:

Note that $\mathbf{s} = [s_0, \dots, s_{N-1}]^T$ specifies one of 2^N basis hypotheses.

Want to find a (small) subset $\mathcal{S}_\star \subset \{0, 1\}^N$ of basis hypotheses with *non-negligible* probability $p(\mathbf{s}|\mathbf{y})$.

2. Signal Estimation:

The MMSE estimate,

$$\hat{\mathbf{x}}_{\text{mmse}} = \sum_{\mathbf{s} \in \{0,1\}^N} p(\mathbf{s}|\mathbf{y}) \underbrace{\mathbb{E}\{\mathbf{x}|\mathbf{y}, \mathbf{s}\}}_{\hat{\mathbf{x}}_{\text{mmse}}|\mathbf{s}} \approx \sum_{\mathbf{s} \in \mathcal{S}_\star} p(\mathbf{s}|\mathbf{y}) \mathbb{E}\{\mathbf{x}|\mathbf{y}, \mathbf{s}\},$$

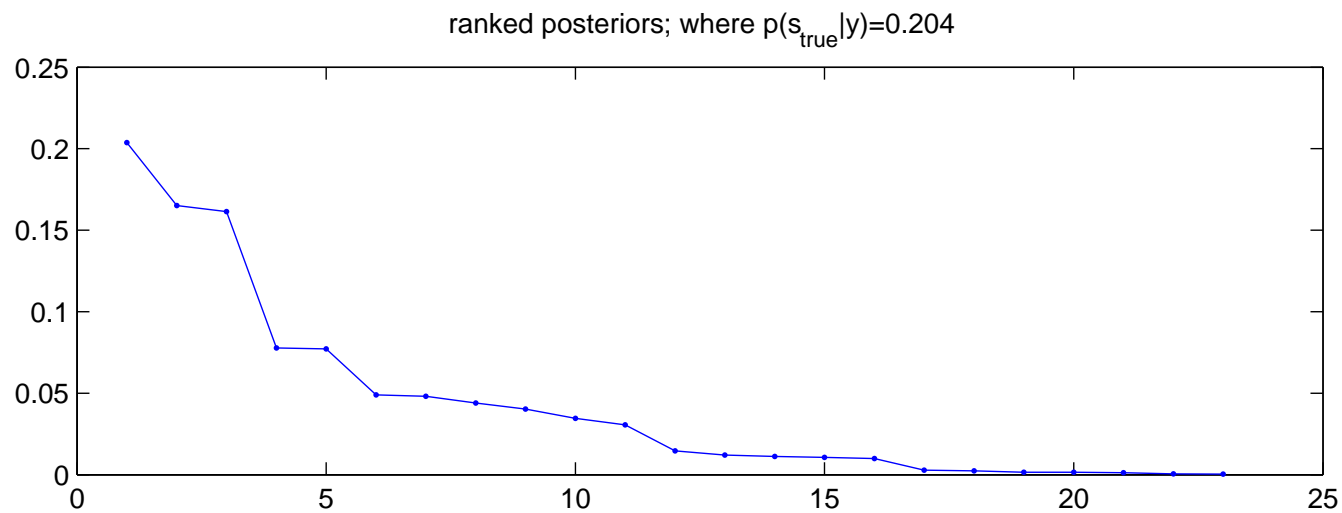
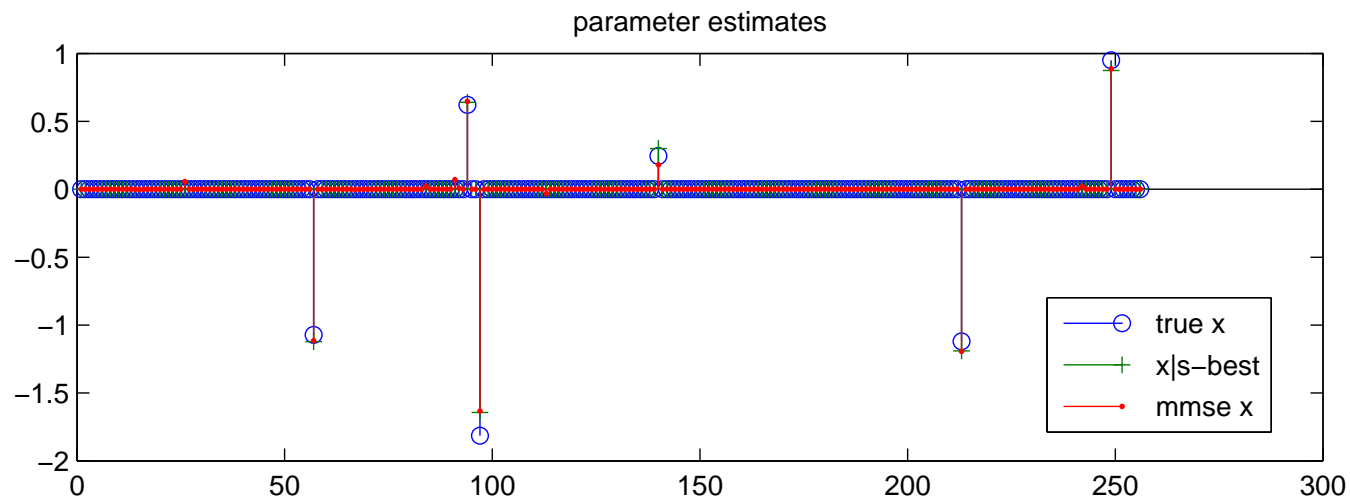
leverages the *inherent uncertainty* in basis estimation.

Can contrast $\hat{\mathbf{x}}_{\text{mmse}}$ with the MMSE estimate conditioned on the most probable basis ($\hat{\mathbf{s}}_{\text{map}}$):

$$\hat{\mathbf{x}}_{\text{mmse}}|\hat{\mathbf{s}}_{\text{map}} = \mathbb{E}\{\mathbf{x}|\mathbf{y}, \hat{\mathbf{s}}_{\text{map}}\}.$$

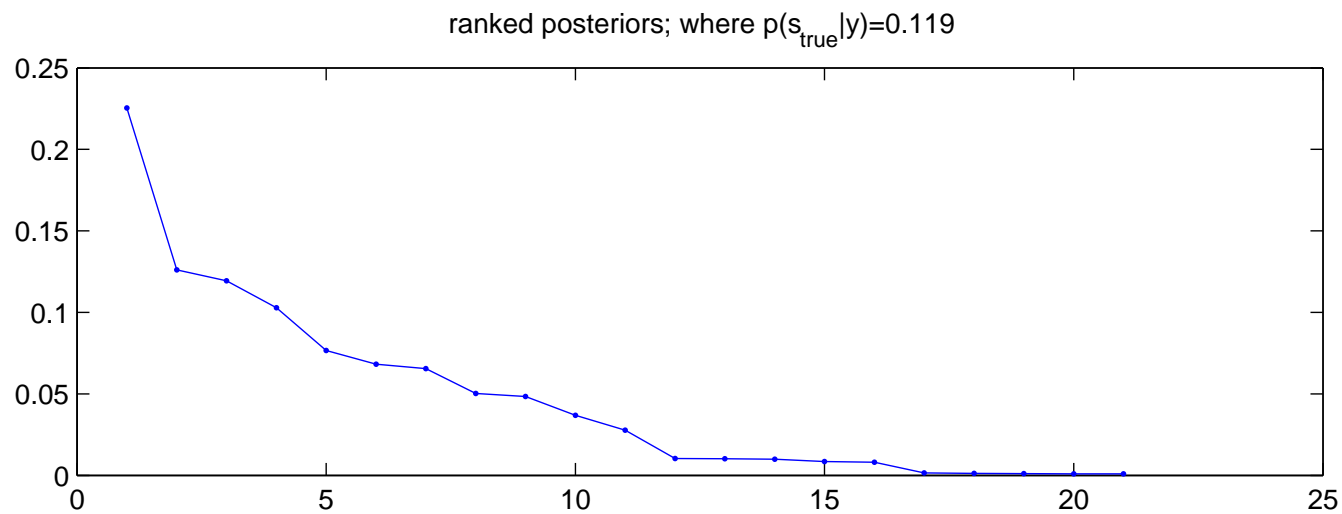
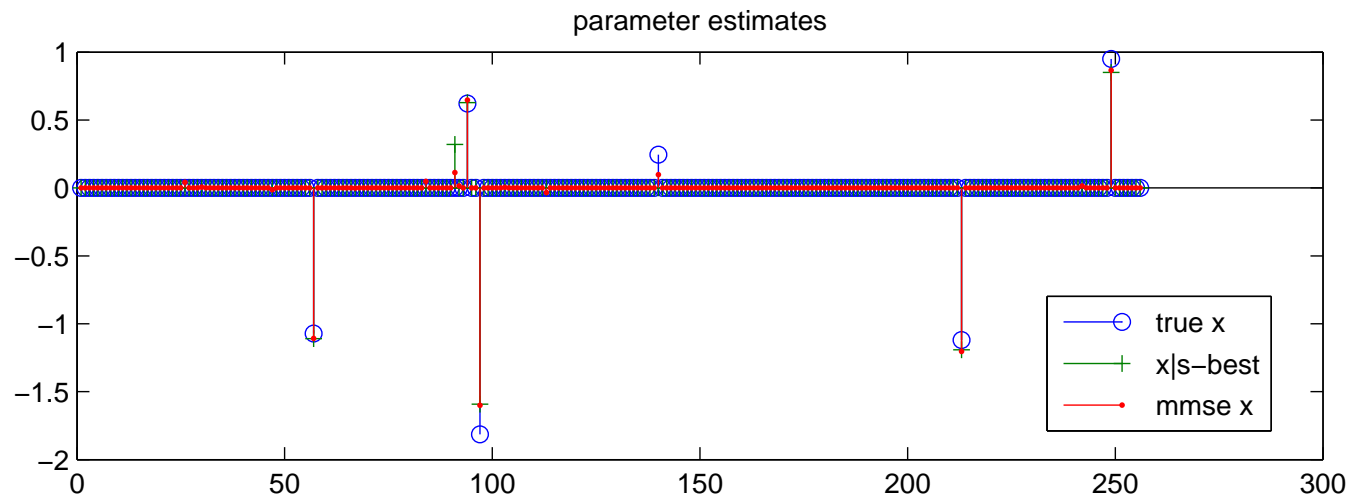
Example – Leveraging Basis Uncertainty:

At SNR=14dB, $\hat{\mathbf{x}}_{\text{mmse}} | \hat{\mathbf{s}}_{\text{map}} \approx \hat{\mathbf{x}}_{\text{mmse}} \approx \mathbf{x}$:



Example – Leveraging Basis Uncertainty (cont.):

But at SNR=13dB, $\hat{\mathbf{x}}_{\text{mmse}} | \hat{\mathbf{s}}_{\text{map}} \neq \hat{\mathbf{x}}_{\text{mmse}} \approx \mathbf{x}!!$



A Basis Selection Metric:

From Bayes rule, we know

$$p(\mathbf{s}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{s})p(\mathbf{s})}{p(\mathbf{y})}.$$

The “basis selection metric”

$$\mu(\mathbf{s}) := \log p(\mathbf{y}|\mathbf{s})p(\mathbf{s}),$$

can be expanded as

$$\begin{aligned} \mu(\mathbf{s}) = & -\frac{M}{2} \ln 2\pi - \frac{1}{2} \ln \det(\boldsymbol{\Sigma}(\mathbf{s})) - \frac{1}{2} \mathbf{y}^T \boldsymbol{\Sigma}(\mathbf{s})^{-1} \mathbf{y} \\ & + \|\mathbf{s}\|_0 \ln \frac{p_1}{1-p_1} + N \ln(1-p_1), \end{aligned}$$

using $\boldsymbol{\Sigma}(\mathbf{s}) := \text{Cov}\{\mathbf{y}|\mathbf{s}\} = \mathbf{A} \mathcal{D}(\mathbf{s}) \mathbf{A}^T + \sigma^2 \mathbf{I}_M$.

So how can we find the basis hypotheses $\{\mathbf{s}\}$ with large $\mu(\mathbf{s})$?

Simple Bayesian Matching Pursuit (BMP-1):

Consider matching pursuit (MP), but using metric $\mu(\mathbf{s})$ rather than magnitude of projection-onto-residual:

$$\mathbf{s}_\star^{(0)} = \mathbf{0};$$

for $i = 1 : P$,

$$\mathcal{S}^{(i)} = \{\text{extensions of } \mathbf{s}_\star^{(i-1)} \text{ to one more active element}\};$$

$$\mathbf{s}_\star^{(i)} = \operatorname{argmax}_{\mathbf{s} \in \mathcal{S}^{(i)}} \mu(\mathbf{s});$$

end;

$$\hat{\mathcal{S}}_\star = \{\mathbf{s}_\star^{(i)}\}_{i=0}^P;$$

$$\hat{\mathbf{x}} = \sum_{\mathbf{s} \in \hat{\mathcal{S}}_\star} \hat{p}(\mathbf{s}|\mathbf{y}) \mathbb{E}\{\mathbf{x}|\mathbf{y}, \mathbf{s}\} \quad \text{for} \quad \hat{p}(\mathbf{s}|\mathbf{y}) = \frac{\exp\{\mu(\mathbf{s})\}}{\sum_{\mathbf{s}' \in \hat{\mathcal{S}}_\star} \exp\{\mu(\mathbf{s}')\}}$$

Can choose P via Gaussian approx of $\|\mathbf{s}\|_0 \sim \text{Binomial}(N, p_1)$:

$$P = \lceil \operatorname{erfc}^{-1}(2\mathcal{P}_0) \sqrt{2Np_1(1-p_1)} + Np_1 \rceil \quad \text{for} \quad \mathcal{P}_0 := \Pr\{\|\mathbf{s}\|_0 > P\}$$

Reminiscent of the technique proposed in

E. Larsson and Y. Selén, “Linear regression with a sparse parameter vector,”
IEEE TSP, Feb. 2007.

Bayesian Matching Pursuit (BMP- D):

Extension of simple BMP to D simultaneous hypotheses:

$$\mathcal{S}_\star^{(0)} = \{\mathbf{0}\};$$

for $i = 1 : P$,

$$\mathcal{S}^{(i)} = \{\text{unique 1-tap extensions of } \mathcal{S}_\star^{(i-1)}\};$$

$$\mathcal{S}_\star^{(i)} = \{\text{the } D \text{ elements of } \mathcal{S}^{(i)} \text{ with largest } \mu(\mathbf{s})\};$$

end;

$$\hat{\mathcal{S}}_\star = \bigcup_{i=0}^P \mathcal{S}_\star^{(i)};$$

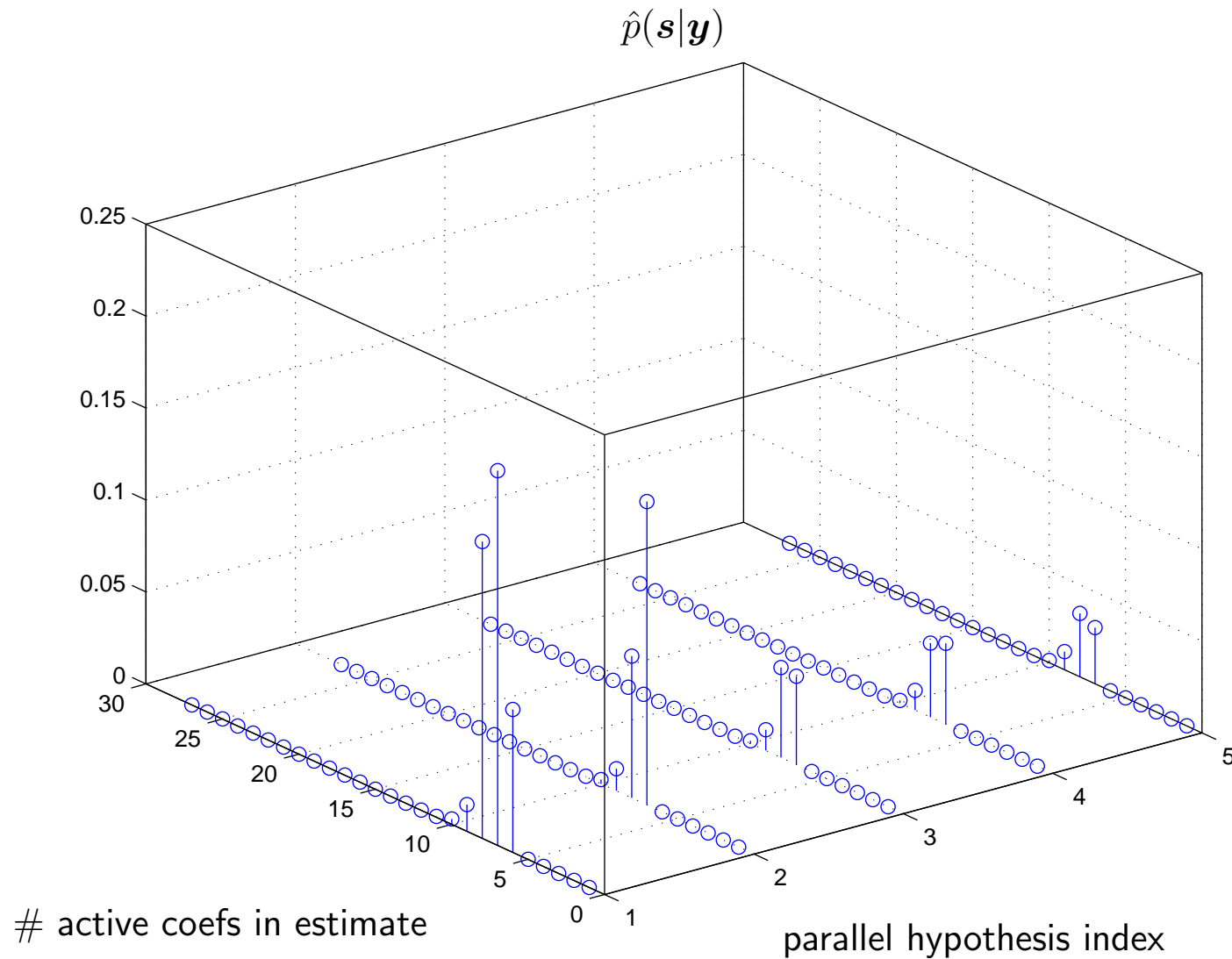
$$\hat{\mathbf{x}} = \sum_{\mathbf{s} \in \hat{\mathcal{S}}_\star} \hat{p}(\mathbf{s}|\mathbf{y}) \mathbb{E}\{\mathbf{x}|\mathbf{y}, \mathbf{s}\} \quad \text{for} \quad \hat{p}(\mathbf{s}|\mathbf{y}) = \frac{\exp\{\mu(\mathbf{s})\}}{\sum_{\mathbf{s}' \in \hat{\mathcal{S}}_\star} \exp\{\mu(\mathbf{s}')\}}$$

D effects a tradeoff between search accuracy and complexity.

A similar extension of basic MP was proposed in

S. F. Cotter and B. D. Rao, "Application of tree-based searches to matching pursuit," *Proc. IEEE ICASSP*, 2001.

Example – Basis-hypothesis posteriors found by BMP- D :



Fast Bayesian Matching Pursuit (FBMP):

- The BMP-like approach proposed by Larsson and Selén consumes $\mathcal{O}(N^3 M^2)$ multiplications.
 \rightsquigarrow Too expensive!!
- Using a recursive metric update, we propose a *fast* BMP-1 that consumes only $\mathcal{O}(NMP)$ multiplications.
 \rightsquigarrow Savings of $\mathcal{O}(N^2 \frac{P}{M})$; many orders of magnitude!
- Can straightforwardly extend to a *fast* BMP- D with complexity $\mathcal{O}(NMPD)$.

Fast Recursive Evaluation of Metric $\mu(\mathbf{s})$:

FBMP's fast update is based on two key properties:

1. $\Delta_n(\mathbf{s})$, the *change* in metric $\mu(\mathbf{s})$ that results from activating the n^{th} tap in \mathbf{s} , can be expressed as

$$\Delta_n(\mathbf{s}) = \frac{1}{2} \log \beta_n^{(i)} + \frac{1}{2} \beta_n^{(i)} (\mathbf{y}^T \mathbf{b}_n^{(i)})^2 + \log \frac{p_1}{1-p_1}$$

$$\beta_n^{(i)} = (1 + \mathbf{a}_n^T \mathbf{b}_n^{(i)})^{-1}$$

$$\mathbf{b}_n^{(i)} = \Sigma(\mathbf{s})^{-1} \mathbf{a}_n \quad \longleftarrow \text{but still } \mathcal{O}(M^2)?$$

2. The structure of $\Sigma(\mathbf{s})^{-1}$ permits $\mathcal{O}(M)$ calculation of $\mathbf{b}_n^{(i)}$:

$$\mathbf{b}_n^{(i)} = \mathbf{b}_n^{(i-1)} - \mathbf{b}_{n_\star}^{(i-1)} \beta_{n_\star}^{(i-1)} \mathbf{b}_{n_\star}^{(i-1)T} \mathbf{a}_n,$$

for n_\star the index of the $\mu(\mathbf{s})$ -maximizing element of $\mathcal{S}^{(i-1)}$. Here, $i := \|\mathbf{s}\|_0$, and so $^{(i-1)}$ refers to the parent node.

FBMP-1:

$$\mu^{(0)} = N \log(1 - p_1) - \frac{M}{2} \log 2\pi - \frac{M}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \|\mathbf{y}\|_2^2;$$

$$\mathcal{N}^{(0)} = \{\};$$

$$\mathbf{b}_n^{(1)} = \frac{1}{\sigma^2} \mathbf{a}_n \text{ for all } n \notin \mathcal{N}^{(0)};$$

for $i = 1 : P$,

$$\beta_n = (1 + \mathbf{a}_n^T \mathbf{b}_n^{(i)})^{-1} \text{ for all } n \notin \mathcal{N}^{(i-1)};$$

$$\mu_n = \mu^{(i-1)} + \frac{1}{2} \log \beta_n + \frac{1}{2} \beta_n (\mathbf{y}^T \mathbf{b}_n^{(i)})^2 + \log \frac{p_1}{1-p_1} \text{ for all } n \notin \mathcal{N}^{(i-1)};$$

$$n_\star^{(i)} = \operatorname{argmax}_n \mu_n;$$

$$\mu^{(i)} = \mu_{n_\star^{(i)}};$$

$$\mathcal{N}^{(i)} = \mathcal{N}^{(i-1)} \cup \{n_\star^{(i)}\};$$

$$\mathbf{b}_n^{(i+1)} = \mathbf{b}_n^{(i)} - \mathbf{b}_{n_\star^{(i)}}^{(i)} \beta_{n_\star^{(i)}} \mathbf{b}_{n_\star^{(i)}}^{(i)T} \mathbf{a}_n \text{ for all } n \notin \mathcal{N}^{(i-1)};$$

end;

$$\hat{\mathbf{x}}_i = \sum_{j=1}^i \delta_{n_\star^{(j)}} \mathbf{b}_{n_\star^{(j)}}^{(i+1)T} \mathbf{y} \text{ for all } i \in \{1, \dots, P\};$$

$$\hat{p}_i = \frac{\exp\{\mu^{(i)}\}}{\sum_{j=0}^P \exp\{\mu^{(j)}\}} \text{ for all } i \in \{1, \dots, P\};$$

$$\hat{\mathbf{x}} = \sum_{i=1}^P \hat{p}_i \hat{\mathbf{x}}_i;$$

Numerical Experiments – Comparison to Other Algs:

Nominal Params: $N = 512$

$p_1 = 0.04$... so $p_1 N = 20$ active coefs on average

$M = 120$

SNR = 19 dB

Note: Considering $M \gtrsim p_1 N \log\left(\frac{N}{p_1 N}\right) \Leftrightarrow \frac{\# \text{ active coefs}}{\text{measurement}} = \frac{p_1 N}{M} \lesssim \frac{-1}{\log(p_1)}$,
our nominal parameters yield $\frac{p_1 N}{M} = 0.16$ and $\frac{-1}{\log(p_1)} = 0.31$.

Algorithms:

SparseBayes - Wipf & Rao

OMP - Tropp & Gilbert

StOMP - Donoho, Tsaig, Drori & Starck

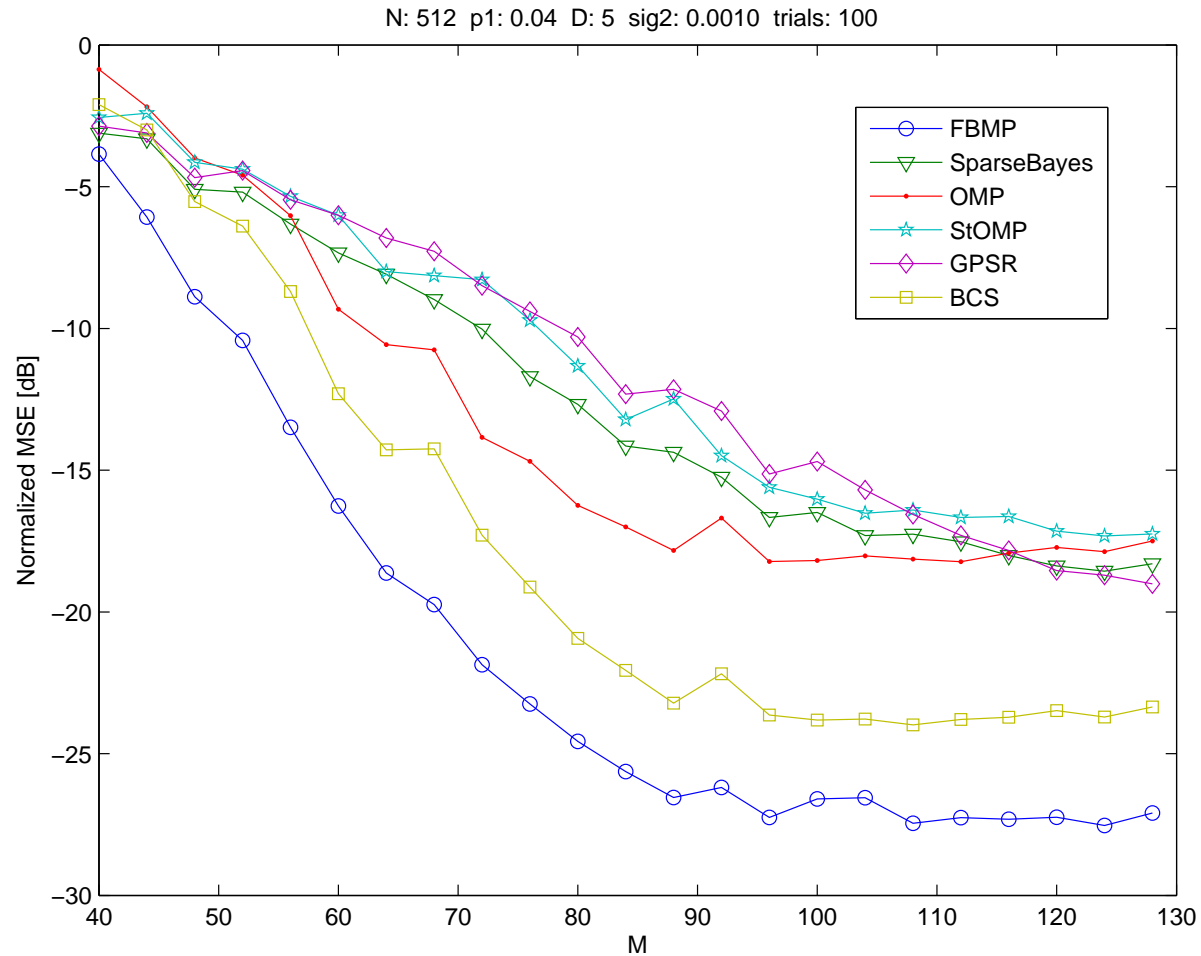
GPSR-Basic - Figueiredo, Nowak & Wright

BCS - Ji & Carin

FBMP - ... with $D = 5$

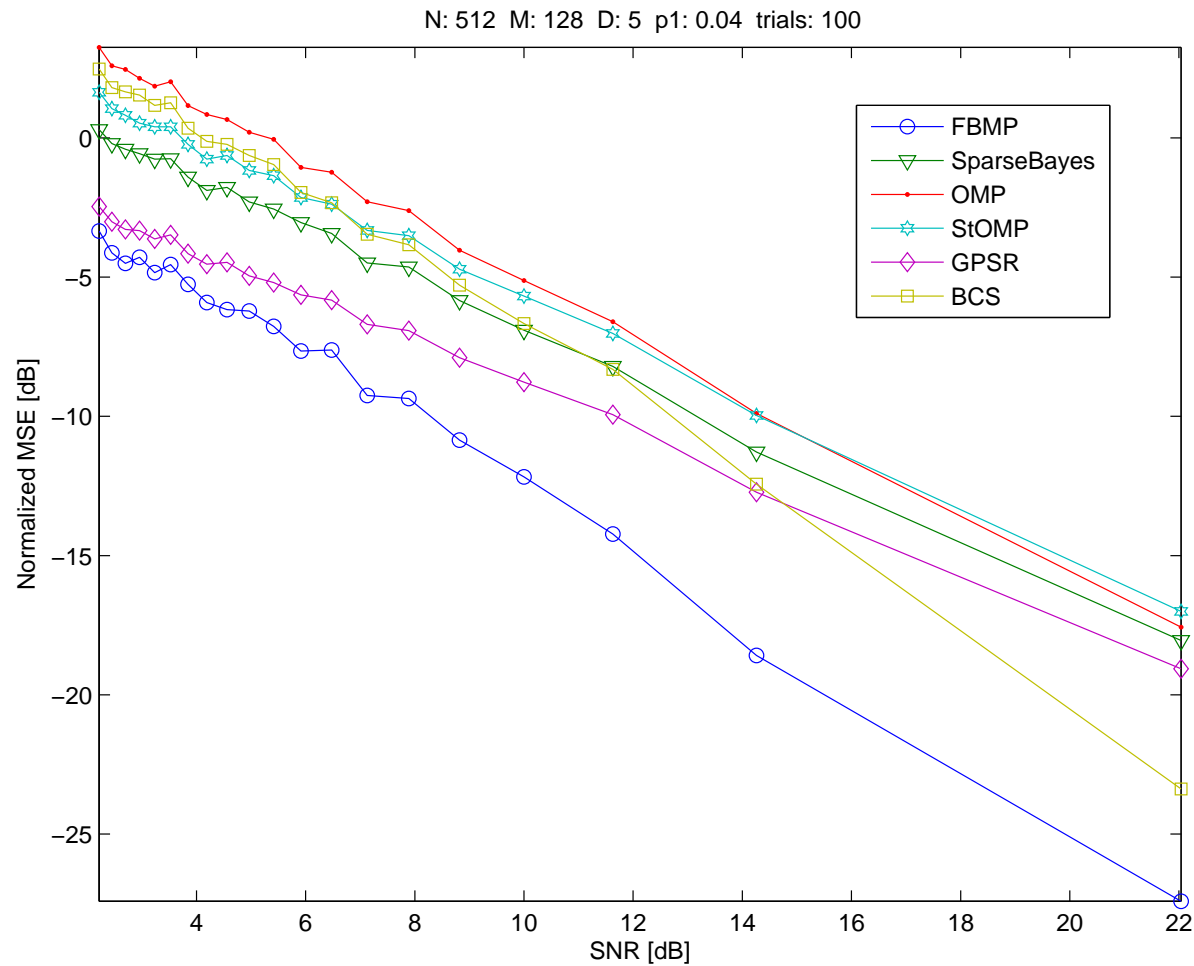
Performance: $\text{NMSE} := \text{Avg} \left\{ \frac{\|\hat{\mathbf{x}} - \mathbf{x}\|_2^2}{\|\mathbf{x}\|_2^2} \right\}$ over 100 random trials.

NMSE versus observation length M :



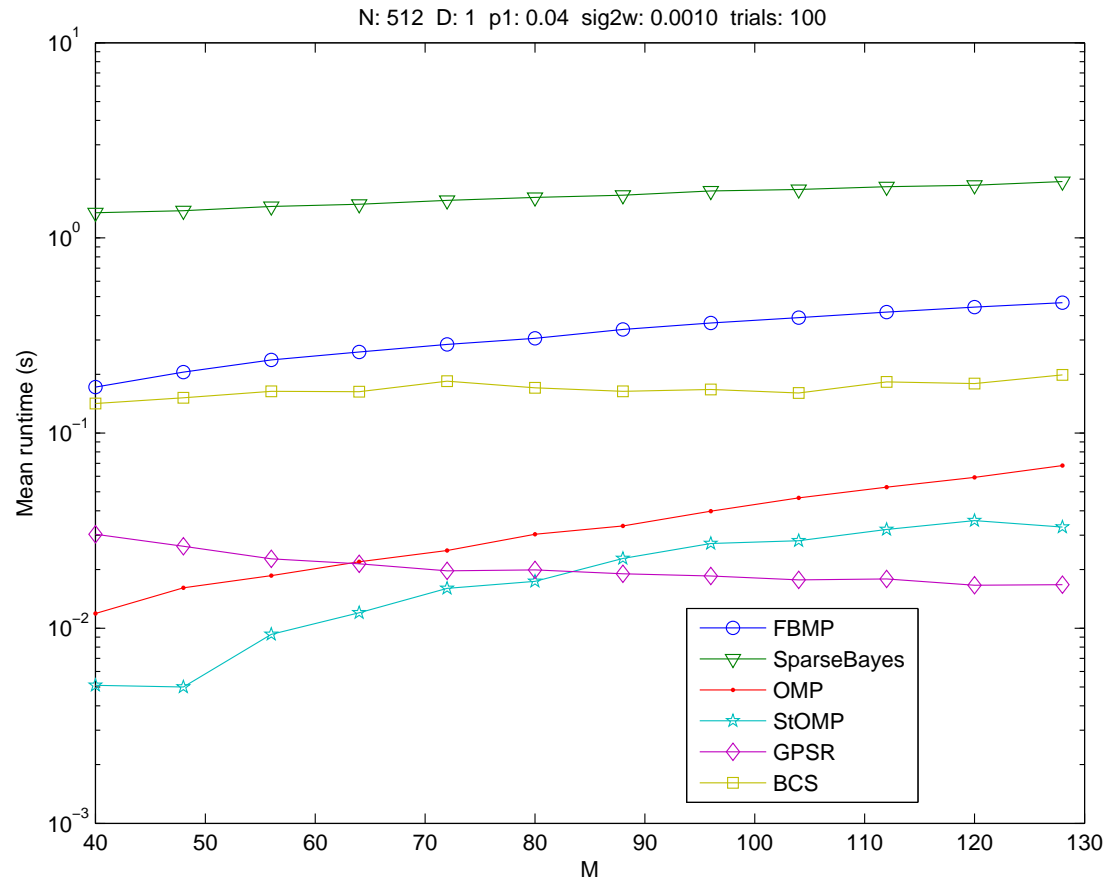
For $\frac{p_1 N}{M} < 0.2$, FBMP outperformed BCS by 3 dB and others by ≥ 10 dB.
 As $\frac{p_1 N}{M} \rightarrow 0.5$, NMSEs converge.

NMSE versus SNR:



At high SNR, FBMP outperformed BCS by 3 dB and others by ≥ 9 dB.
As $\text{SNR} \rightarrow 0$ dB, GPSR catches up.

Runtime versus observation length M :



(Not-yet-optimized) FBMP is an order of magnitude faster than SparseBayes, about the same speed as BCS, and an order of magnitude slower than OMP, StOMP, and GPSR.

Numerical Experiments – FBMP Behavior:

Nominal Signal Parameters:

$$N = 256$$

$$p_1 = 0.04 \quad \dots \text{thus } p_1 N = 10 \text{ active coefs on average}$$

$$M = 64$$

$$\text{SNR} = 15 \text{ dB} \quad \dots \text{where } \text{SNR} := \frac{p_1 N}{\sigma^2 M}$$

\mathbf{A} : i.i.d. $\mathcal{N}(0, 1)$, then columns normalized.

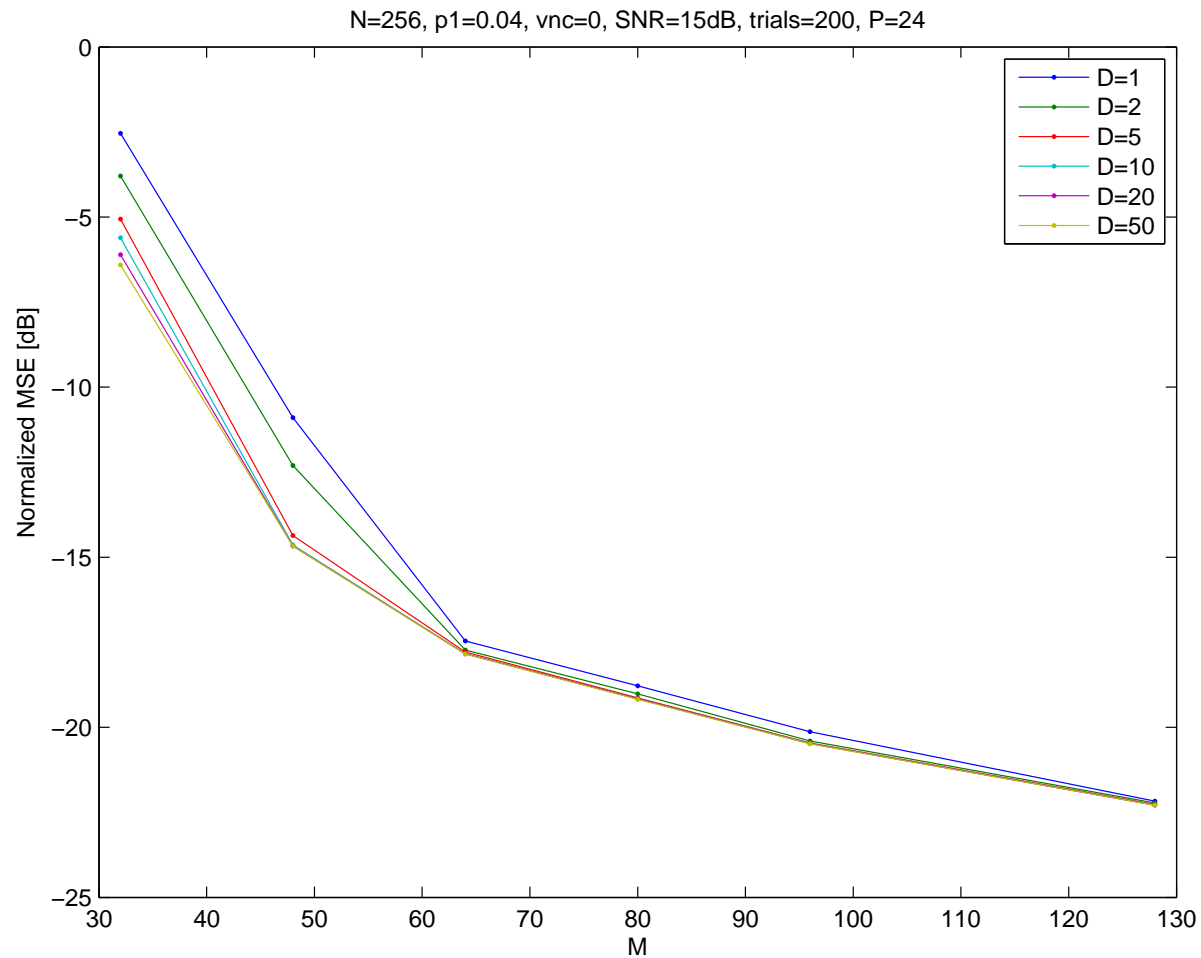
Note:

Considering $M \gtrsim p_1 N \log\left(\frac{N}{p_1 N}\right) \Leftrightarrow \frac{\# \text{ active coefs}}{\text{measurement}} = \frac{p_1 N}{M} \lesssim \frac{-1}{\log(p_1)}$,
 our nominal parameters yield $\frac{p_1 N}{M} = 0.16$ and $\frac{-1}{\log(p_1)} = 0.31$.

Performance:

$$\text{NMSE} := \text{Avg} \left\{ \frac{\|\hat{\mathbf{x}} - \mathbf{x}\|_2^2}{\|\mathbf{x}\|_2^2} \right\} \text{ over 200 random trials.}$$

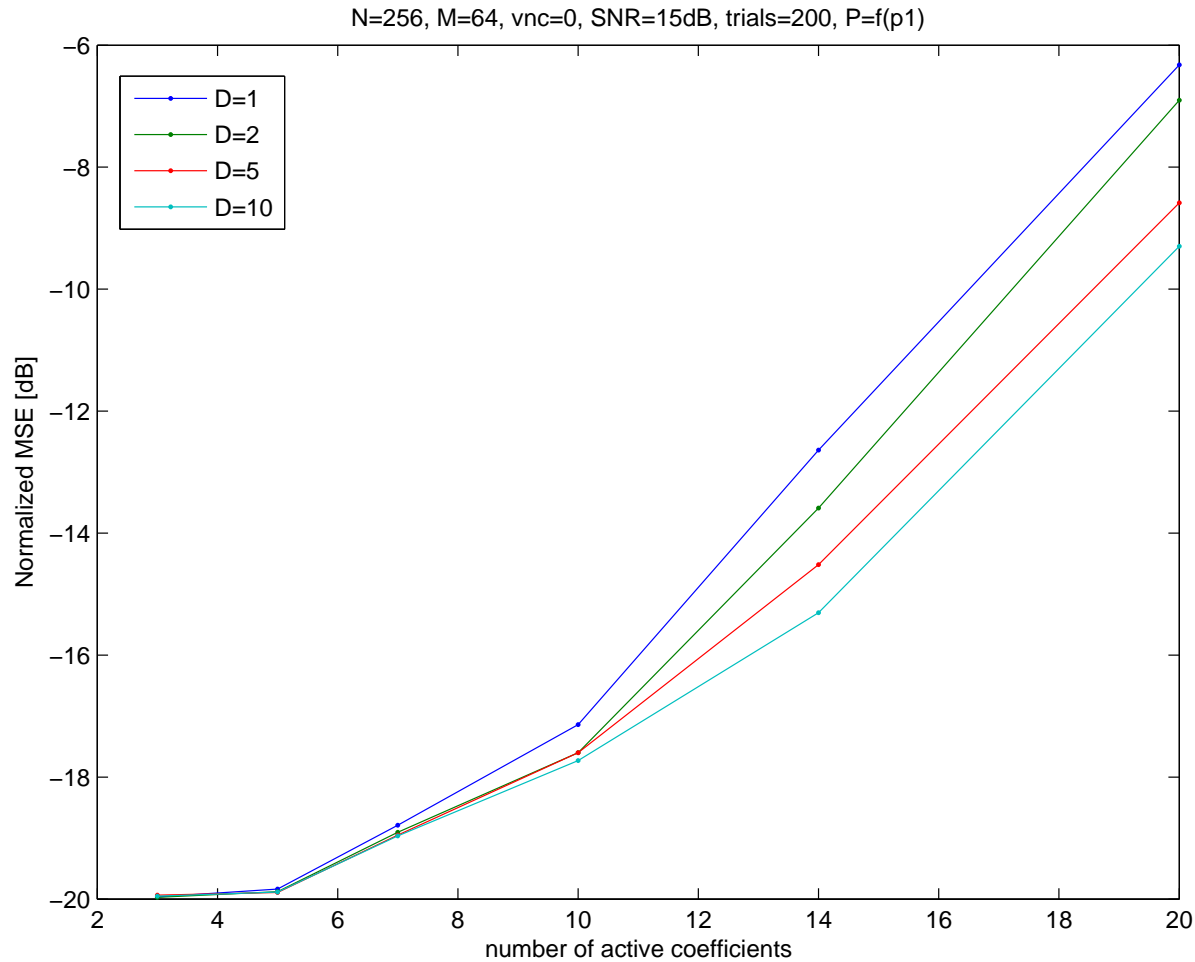
NMSE versus observation length M :



When $D = 1$, knee in curve at $\frac{p_1 N}{M} = \frac{10}{64} = 0.16$ $\frac{\# \text{ active coefs}}{\text{measurement}}$.

For larger D , knee moves to 0.2 $\frac{\# \text{ active coefs}}{\text{measurement}}$ and NMSE improves by 3 dB.

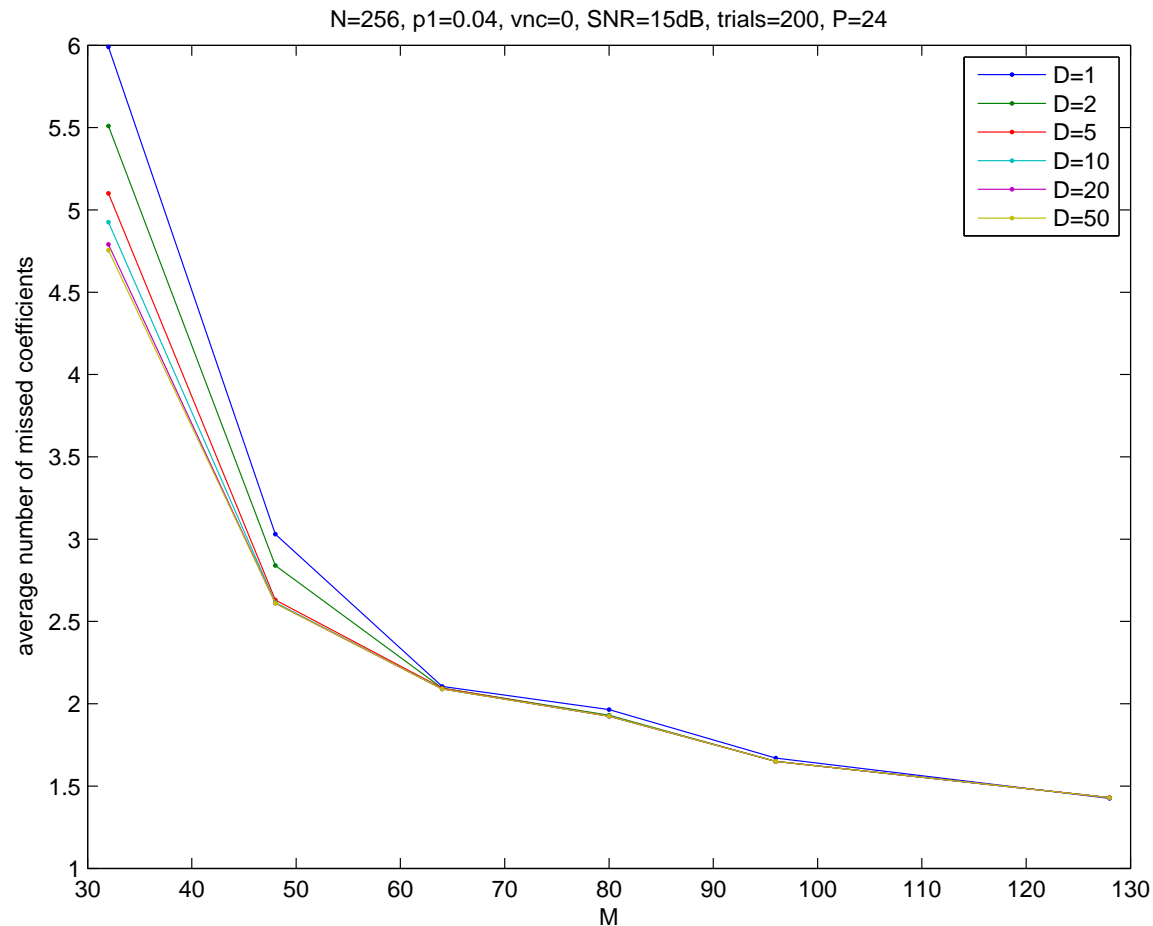
NMSE versus expected # active coefs $p_1 N$:



When $D = 1$, knee in curve at $\frac{p_1 N}{M} = \frac{10}{64} = 0.16 \frac{\# \text{ active coefs}}{\text{measurement}}$.

For $D = 10$, knee at $\frac{14}{64} = 0.2 \frac{\# \text{ active coefs}}{\text{measurement}}$ and NMSE 3 dB improved.

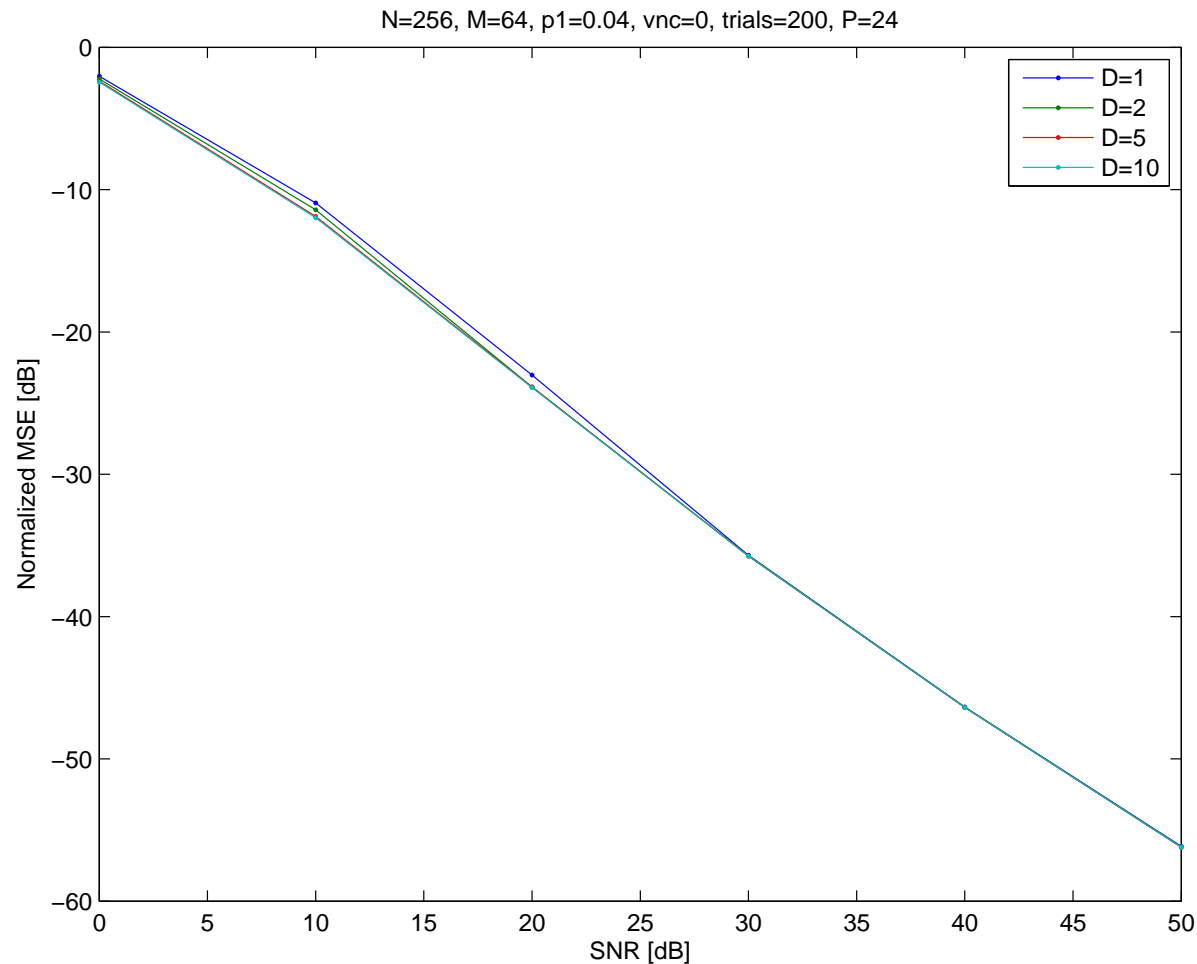
Active coefs missing from \hat{s}_{map} :



Again, knee in curve at $\frac{p_1 N}{M} \approx 0.2 \frac{\# \text{ active coefs}}{\text{measurement}}$.

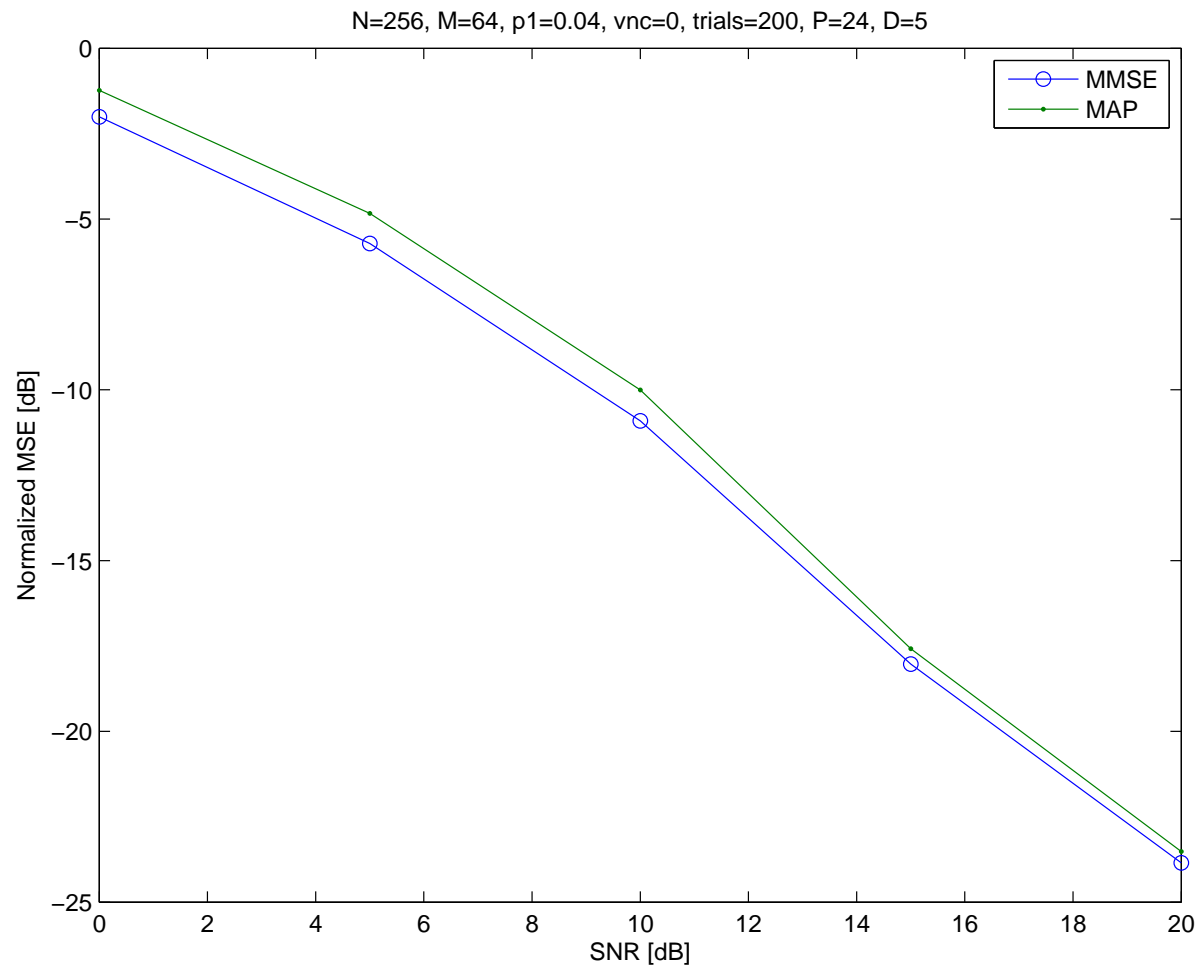
(Note: we expect some misses since signal model is zero-mean.)

NMSE versus SNR:



Note linear dependence between NMSE [dB] & SNR [dB].
(No benefit from D -increase expected since $\frac{p_1 N}{M} = 0.16$.)

NMSE for $\hat{\mathbf{x}}_{\text{mmse}}$ and $\hat{\mathbf{x}}_{\text{mmse}} | \hat{\mathbf{s}}_{\text{map}}$:



Exploiting basis *uncertainty* gives ≈ 1 dB gain in NMSE.

Conclusion:

- Building on the Bayesian sparse-coefficient estimation technique of Larsson and Selén, we proposed
 - a forward search, with fast recursive update, that reduces complexity from $\mathcal{O}(N^3M^2)$ to $\mathcal{O}(NMP)$, and
 - an extension of the search to $D > 1$ simultaneous hypotheses, providing up to 3 dB of NMSE gain when $\frac{\# \text{ active coefs}}{\text{measurement}}$ is high.
- Comparisons against SparseBayes, BCS, OMP, StOMP, GPSR showed NMSE improvements of several dB over a wide range of parameters.
- In the near future, we plan to optimize FBMP and extend our approach to signals with complex-valued coefs modeled with *non-zero means*.