

Statistical Signal Processing

Lecture 12

Statistical Signal Modeling, Learning and Processing

chapter 6: Learning Topics

- Expectation-Maximization (EM) algorithm
- compressed sensing and sparsity
- compressed sensing algorithms: LASSO etc
- Sparse Bayesian Learning (SBL), Empirical Bayes
- Variational Bayes (VB)

Outline

1 Parametric Signal Models

2 Compressed Sensing

3 SBL

Main Messages

- There are many Bayesian estimation problems, many of which are LMMSE (Wiener, Kalman), which contain hyperparameters to be tuned, using various approaches.
- Information combining: from weighted least-squares to message passing in a more general overall Bayesian formulation (e.g. cooperative location estimation)
- Empirical Bayes (EB) as principled framework for bias-variance trade-off
- but not necessarily using empirical Bayes for hyperparameter estimation: SURE, Cross Validation
- compressive sensing, sparse models, generalization of model order selection to support region, model complexity and structure
- Sparse Bayesian Learning (SBL) is one EB instance, allowing to exploit (approximate) sparsity for compressed sensing
 - can be extended to time-varying scenarios with sparse variations also
 - can be extended to dictionary learning, in particular with Kronecker structured dictionaries
- message passing (approximate iterative) inference techniques: easy to get the mean (estimate) correct but more difficult to get correct posterior variances

Kalman Filter

Linear state-space model:

state update equation:

$$\mathbf{x}_{k+1} = \mathbf{F}_k(\theta) \mathbf{x}_k + \mathbf{G}_k(\theta) \mathbf{w}_k$$

measurement equation:

$$\mathbf{y}_k = \mathbf{H}_k(\theta) \mathbf{x}_k + \mathbf{v}_k$$

for $k = 1, 2, \dots$, with uncorrelated

- initial state $\mathbf{x}_0 \sim \mathcal{N}(\hat{\mathbf{x}}_0, \mathbf{P}_0)$,
- measurement noise $\mathbf{v}_k \sim \mathcal{N}(0, \mathbf{R}_k(\theta))$,
- state noise $\mathbf{w}_k \sim \mathcal{N}(0, \mathbf{Q}_k(\theta))$.

State model known up to some parameters θ .

Often $\mathbf{F}_k(\theta)$, $\mathbf{G}_k(\theta)$, $\mathbf{H}_k(\theta)$ linear in θ : bilinear case.



Numerous Applications

- LMMSE wireless channel estimation:

x_k = FIR filter response, θ : Power Delay Profile, AR(1) dynamics in e.g. diagonal F and Q

- Bayesian adaptive filtering (BAF):

analogous to LMMSE channel estimation, except measurement equation: only one 1D measurement is available per instance. An extremely simplified form of BAF is the so-called **Proportionate LMS (P-LMS)** algorithm.

- Position tracking (GPS):

$$\mathbf{x}_{t+1} = \begin{bmatrix} 1 & \Delta t & \frac{1}{2}\Delta t^2 \\ 0 & 1 & \Delta t \\ 0 & 0 & 1 \end{bmatrix} \cdot \mathbf{x}_t = \begin{bmatrix} x_t + \Delta t \cdot v_t + \frac{1}{2}\Delta t^2 a \\ v_t + \Delta t \cdot a \\ a \end{bmatrix}$$

θ : acceleration model parameters (e.g. white noise, AR(1))

- Blind Audio Source Separation (BASS):

x_k = source signals,

θ : (short+long term) AR parameters, reverb filters



Static LMMSE (Wiener) Applications

- Direction of Arrival (DoA) estimation: x = decorrelated sources, apart from the DoA parameters there could also be antenna array calibration parameters or source and noise covariance parameters.
- Blind and semi-blind channel estimation. In the techniques that exploit the (white) second-order statistics of x , (the unknown part of) x gets modeled as Gaussian. Numerous variations: single-carrier, OFDM and CDMA versions, single- and multi-user, single- and multi-stream, with black box FIR channel models or propagation based parameterized channel models.
Image Deblurring, NMRI Imaging
- LMMSE receiver (Rx) design: x = Tx symbol sequence to be recovered on the basis of Rx signal, in single- or multi-user settings and other variations as in the channel estimation case. The crosscorrelation between Tx and Rx signals depends on the channel response, which is part of the parameters. The Rx signal covariance matrix on the other hand can be modeled in various ways, non-parametric or parametric. Account for the channel estimation error in the LMMSE Rx design. Another approach: consider the LMMSE filter directly as the parameters.
LMMSE Tx design, partial CSIR/CSIT.

Adaptive Kalman Filter solutions

- Extended Kalman Filter (**EKF**)
- other generic nonlinear Kalman Filter extensions:
Unscented Kalman Filter (**UKF**), Cubature Kalman Filter (**CKF**), Gaussian Sum Filter, Particle Filter (**PF**)
- Recursive Prediction Error Method (**RPEM**) Kalman Filter
- Second-Order Extended Kalman Filter (**SOEKF**)
- Expectation-Maximization (**EM**)/Variational Bayes (**VB**) Kalman Filter

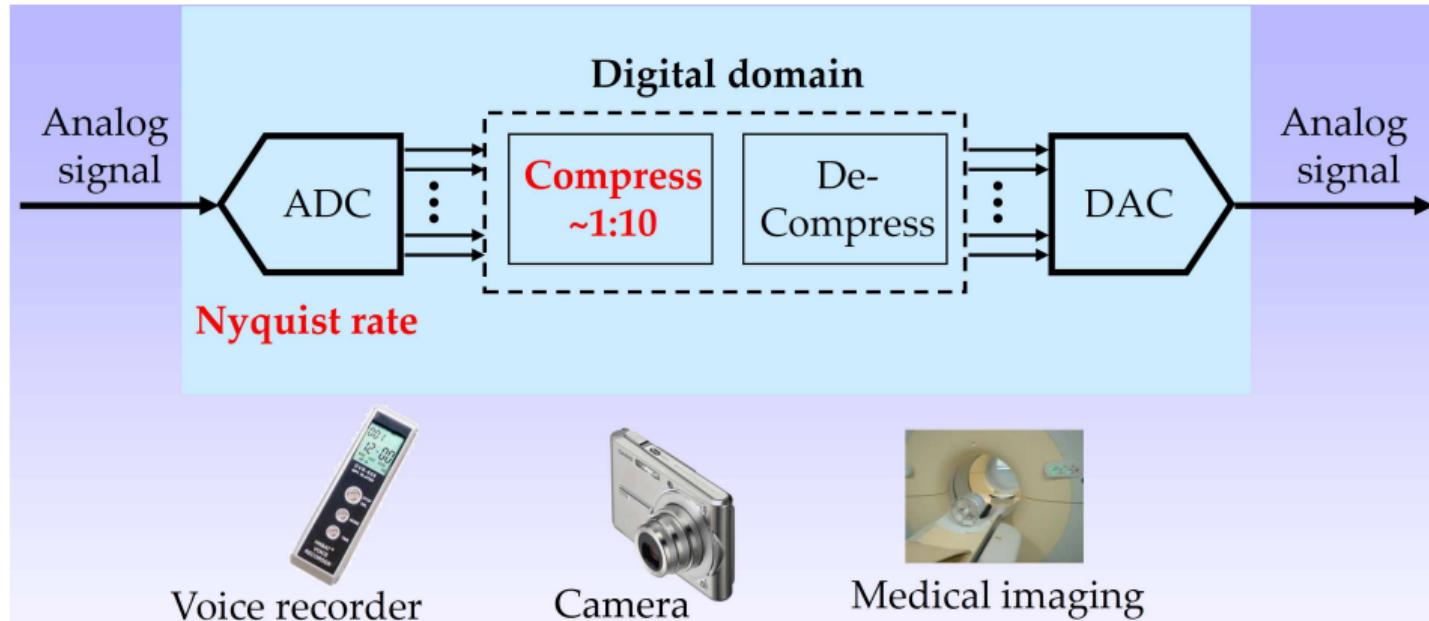
Outline

1 Parametric Signal Models

2 Compressed Sensing

3 SBL

Compressed Sensing



Compressed Sensing

In some applications, measurements are costly:

- Magnetic resonance imaging:
 - scan time \approx 30 minutes
 - scan time proportional to # samples taken

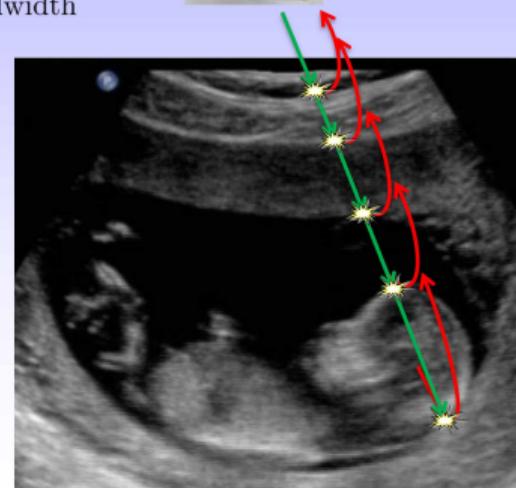
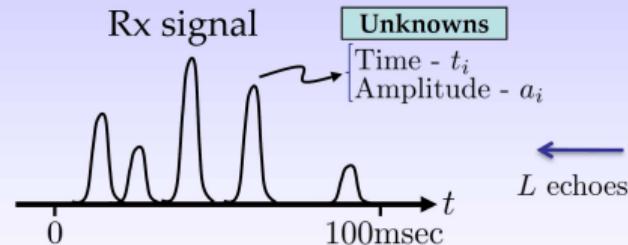


- Imaging outside visible spectrum:
 - CMOS doesn't work
 - high cost per pixel
- Wireless communication:
 - pilots inserted to measure channel
 - more pilots means less payload



Ultra-Sound Imaging

- High sampling rates
- High digital processing rates

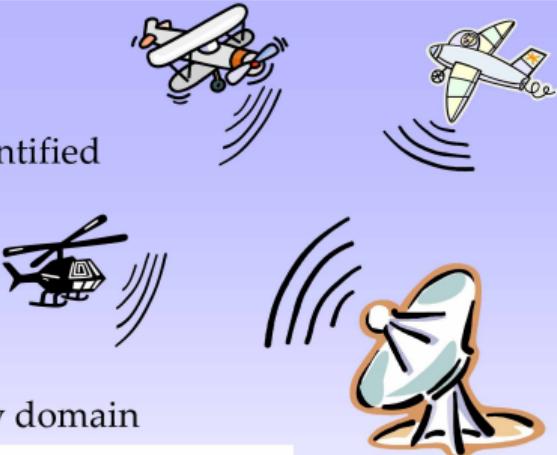


- Echoes result from scattering in the tissue
- The image is formed by identifying the scatterers

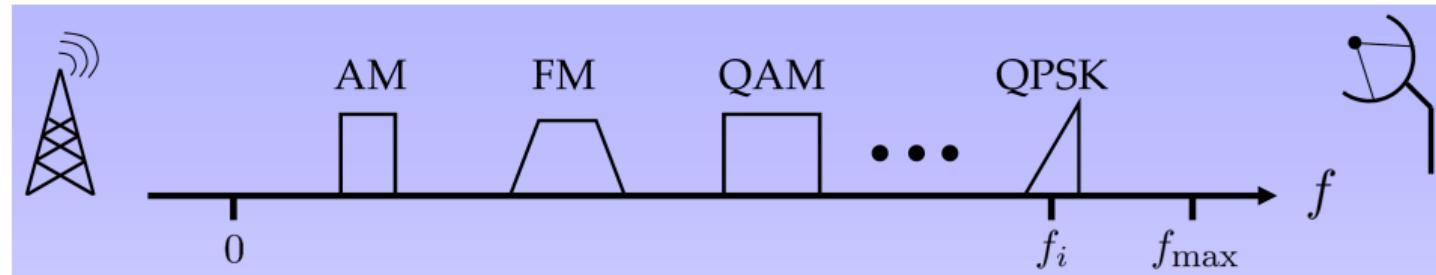
Radar

- Principle:
 - A known pulse is transmitted
 - Reflections from targets are received
 - Target's ranges and velocities are identified

 - Challenge:
 - All processing is done digitally
 - Targets can lie on an arbitrary grid
 - Process of digitizing
→ loss of resolution in range-velocity domain

 - Subspace methods:
- 
- The diagram illustrates various radar configurations. At the top left, two aircraft are shown with their respective signal waveforms. Below them, a helicopter and a ground-based satellite dish are also depicted with their signal patterns. To the right of the challenge section is a 2D heatmap showing the relationship between Delay ($\times \tau_{\max}$) on the x-axis and Doppler ($\times v_{\max}$) on the y-axis. The plot features several data points: red 'x' marks for 'True Targets' and blue open circles for 'MF peaks'. A color scale bar on the right indicates intensity from 0.2 (blue) to 1.0 (red). The axes range from 0 to 1.0.

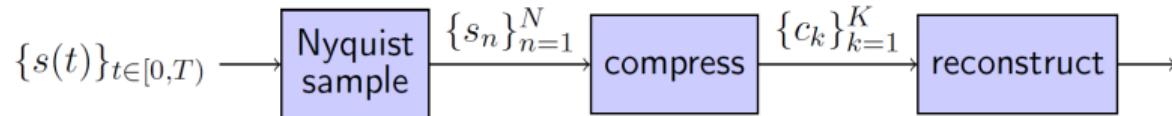
Cognitive Radio



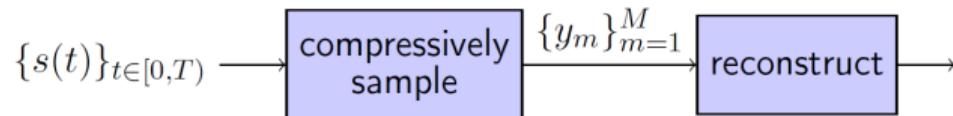
- The spectrum occupation by legacy/primary users is sparse.
- Unlicensed secondary users can insert in spectral holes.
- However, have to find the spectral holes, or the occupied spectrum portions.
- Generalized Nyquist says that one can sample at a rate exceeding the sum of the bandwidth, however here also (only) the spectral support needs to be estimated.

Compressive vs Classical Sampling

- Classical approach:



- New approach:



$$\text{Nyquist rate } \frac{N}{T} \gg \text{compressive sampling rate } \frac{M}{T} \gtrsim \text{information rate } \frac{K}{T}$$

Linear Measurements Model

- 1 For now, assume noiseless linear measurements, e.g.,

$$y_m = \int_0^T \phi_m(t) s(t) dt, \quad m = 1, \dots, M$$

- 2 Also assume signal $s(t)$ is bandlimited, in which case Nyquist says

$$s(t) = \sum_{n=1}^N s_n \operatorname{sinc}\left(\frac{t}{T_s} - n + 1\right), \quad t \in [0, T).$$

Putting these together, we get the convenient discrete representation

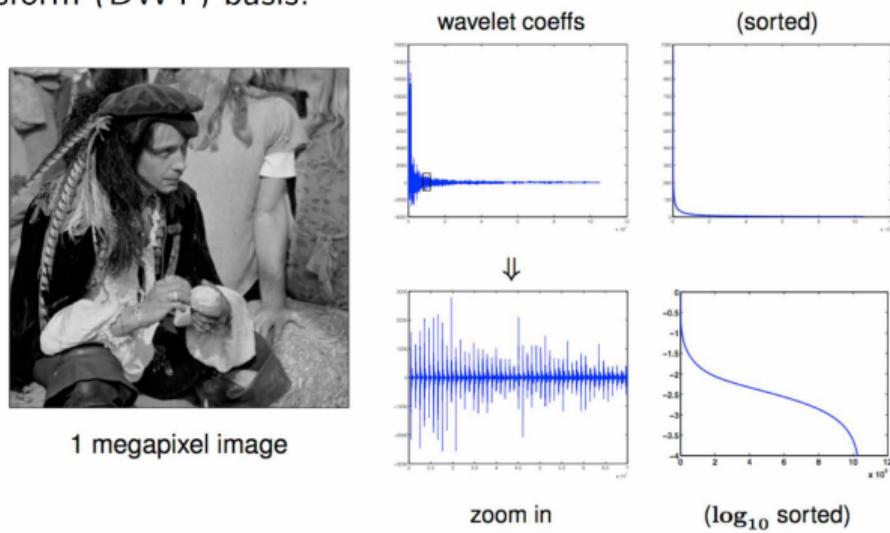
$$y_m = \sum_{n=1}^N s_n \underbrace{\int_0^T \phi_m(t) \operatorname{sinc}\left(\frac{t}{T_s} - n + 1\right) dt}_{\triangleq \Phi_{m,n}}$$

or, in matrix/vector form, $\boxed{\mathbf{y} = \Phi s}$ for $s \in \mathbb{R}^N$ and $\mathbf{y} \in \mathbb{R}^M$.



Sparsity

- Many real-world signals are approximately sparse in a known basis.
- For example, natural images are sparse in the discrete wavelet transform (DWT) basis:



Typically: 99% signal energy captured by only 2.5% of DWT coefficients!

Sparsity will be captured by prior information in a Bayesian approach.

K -sparse in a Dictionary Ψ

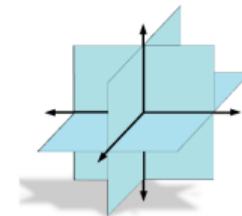
- We say that a signal class \mathcal{S} is K -sparse in the dictionary Ψ if each $s \in \mathcal{S}$ can be written as

$$s = \Psi x$$

for some K -sparse vector x (i.e., x has at most K nonzero elements).

- Usually orthonormal dictionaries Ψ are used (e.g., DWT, DCT, DFT), but overcomplete dictionaries may also be considered.

- Geometrically, a K -sparse vector $x \in \mathbb{R}^N$ lives in a union of $\binom{N}{K}$ subspaces, each of dimension K :



Combining Sparsity with Compressed Measurements

Recall...

- Linear measurement model: $y = \Phi s$ for $\Phi \in \mathbb{R}^{M \times N}$
- Sparse signal model: $s = \Psi x$ for K -sparse $x \in \mathbb{R}^N$

Together...

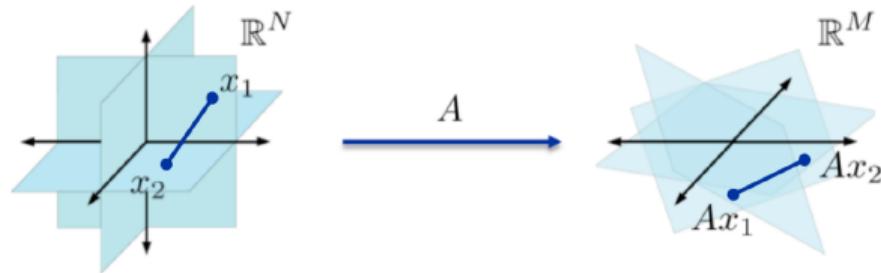
- Compressive sensing model: $y = \underbrace{\Phi \Psi}_{\triangleq A} x$ for $A \in \mathbb{R}^{M \times N}$

Questions:

- 1 What properties of A ensure the recovery of x ?
- 2 Given dictionary Ψ , how can we design Φ to ensure a good A ?

Restricted Isometry Property

- Recall model: $y = Ax$ for $A \in \mathbb{R}^{M \times N}$ and K -sparse $x \in \mathbb{R}^N$.
- Note: if signals $x_1 \neq x_2$ map to the same y , they can't be recovered!



- In general, for our measurement system to be **information preserving**, we want that $\|x_1 - x_2\|_2 \approx \|Ax_1 - Ax_2\|_2$ for all K -sparse x_1, x_2 , or

$$1 - \delta \leq \frac{\|Ad\|_2^2}{\|d\|_2^2} \leq 1 + \delta \quad \text{for all } 2K\text{-sparse } d. \quad \text{"RIP"}$$

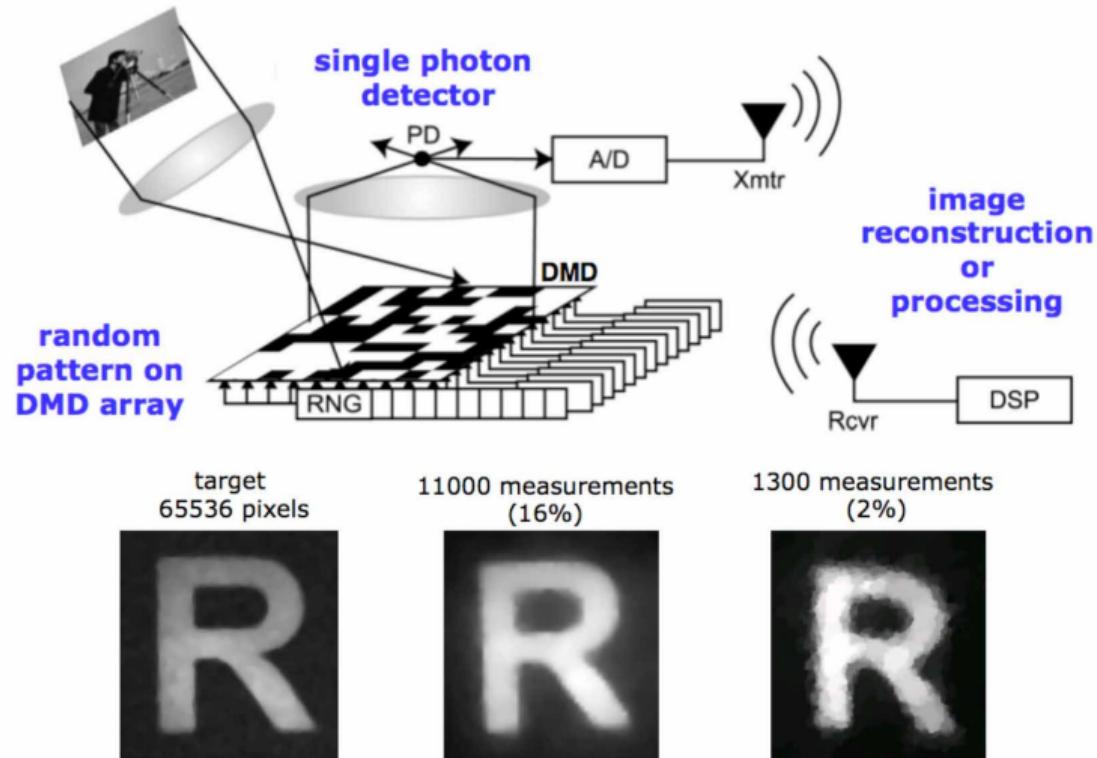
RIP from Randomness

- Testing a given matrix for RIP is an NP-hard (combinatorial) problem.
- Fortunately, if A is **randomly** drawn with **independent zero-mean sub-Gaussian** entries (e.g., normal, ± 1), then *with high probability* it will satisfy RIP if

$$M \geq O\left(K \log \frac{N}{K}\right).$$

- Similarly, if Φ is constructed randomly in the same way, then $A = \Phi\Psi$ will satisfy RIP for any orthonormal Ψ .
- In practice, **semi-random** Φ are preferable, e.g.,
Create $\Phi = JFD$, where D is a diagonal matrix with random ± 1 s,
 F is the N -FFT matrix, and J randomly selects M outputs.

Single Pixel Camera (Rice U.)



Best Sparse Fit - the ℓ_0 technique

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

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

Unfortunately, this is **NP-hard**; we'd need to check all $\binom{N}{K} \approx N^K$ possible supports!

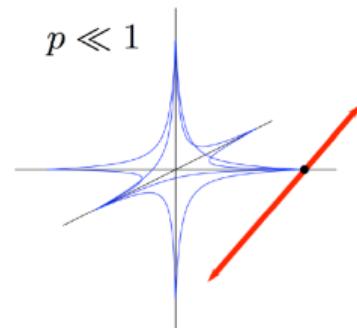
Let's think about this problem geometrically...

Geometry of constrained ℓ_p minimization

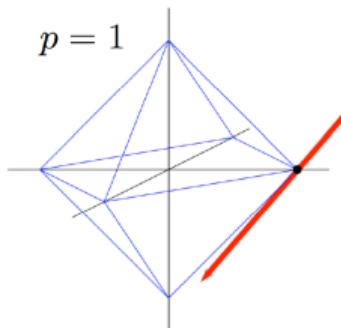
Now consider, for some fixed $p > 0$, the optimization problem:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \underbrace{\|\mathbf{x}\|_p}_{\sqrt[p]{\sum_n |x_n|^p}} \quad \text{s.t. } \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2 \leq \epsilon.$$

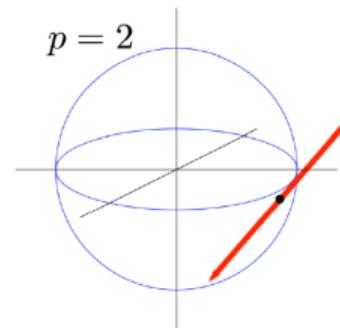
The solution can be found by growing the ℓ_p -ball until it touches the ϵ -rod:



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!

LASSO

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

- Convex! Can be solved very efficiently.
- For \mathbf{A} satisfying $2K$ -RIP, LASSO guarantees that

$$\|\hat{\mathbf{x}} - \mathbf{x}\|_2 \leq \frac{C_1}{\sqrt{K}} \|\mathbf{x} - \mathbf{x}_K\|_1 + C_2 \|\mathbf{w}\|_2$$

where \mathbf{x}_K is the best K -sparse approximation of \mathbf{x} and C_1, C_2 are constants that depend on the RIP δ . Wow!

- In the special case when \mathbf{x} is K -sparse, this simplifies to

$$\|\hat{\mathbf{x}} - \mathbf{x}\|_2 \leq C_2 \|\mathbf{w}\|_2.$$

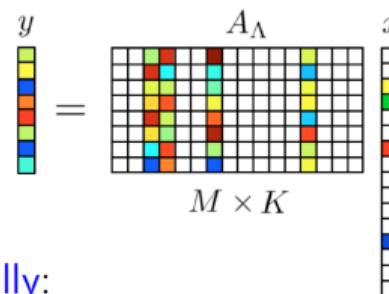
Greedy Search

Main ideas:

- If we can correctly recover the **support** Λ of x (i.e., the locations of nonzeros), then determining the non-zero amplitudes is easy, e.g.,

$$\mathbf{x}_\Lambda = (\mathbf{A}_\Lambda^H \mathbf{A}_\Lambda)^{-1} \mathbf{A}_\Lambda^H \mathbf{y}$$

(least squares)



- Estimate the support **sequentially**:
 - Find the column of \mathbf{A} most “similar” to \mathbf{y} and store its index.
 - Subtract the effect of this column from \mathbf{y} .
 - Repeat (until residual is sufficiently small)!

Famous algorithms include MP, OMP, IHT, CoSaMP, Subspace Pursuit

Bayesian Methods

In the Bayesian approach, one . . .

- models the signal using a **prior** pdf $p(\mathbf{x})$,
- models the measurement process using a **likelihood** function $p(\mathbf{y}|\mathbf{x})$,
- performs inference via **Bayes rule**, yielding the **posterior** pdf

$$p(\mathbf{x}|\mathbf{y}) = Z^{-1} p(\mathbf{y}|\mathbf{x}) p(\mathbf{x}) \text{ where } Z \text{ is a scaling constant,}$$

- often summarizes the posterior pdf by a **point estimate** like

$$\hat{\mathbf{x}} = \int \mathbf{x} p(\mathbf{x}|\mathbf{y}) d\mathbf{x} \quad \text{MMSE estimate}$$

$$\hat{\mathbf{x}} = \arg \max_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}) \quad \text{MAP estimate}$$

and possibly other statistics that quantify **estimate uncertainty**.

LASSO as Bayesian Method

If we assume ...

- additive white Gaussian noise of variance σ^2
- i.i.d Laplacian signal with rate λ/σ^2

then

- likelihood: $p(\mathbf{y}|\mathbf{x}) = \frac{1}{(2\pi\sigma^2)^{M/2}} \exp(-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{Ax}\|_2^2)$
- prior: $p(\mathbf{x}) = \frac{1}{(2\sigma^2/\lambda)^M} \exp(-\frac{\lambda}{\sigma^2} \|\mathbf{x}\|_1)$

for which the maximum a posteriori (MAP) estimate is

$$\begin{aligned}\hat{\mathbf{x}} &= \arg \max_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}) = \arg \max_{\mathbf{x}} \log (Z^{-1} p(\mathbf{y}|\mathbf{x}) p(\mathbf{x})) \\ &= \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{Ax}\|_2^2 + \lambda \|\mathbf{x}\|_1\end{aligned}$$

which is an unconstrained version of the LASSO problem.

Relevance Vector Machine - Sparse Bayesian Learning

- The RVM is based on the *conditionally Gaussian* priors

$$p(\mathbf{x}|\boldsymbol{\alpha}) = \prod_{n=1}^N \mathcal{N}(x_n; 0, \alpha_n^{-1}) \quad \text{and} \quad p(\boldsymbol{\alpha}) = \prod_{n=1}^N \Gamma(\alpha_n; 0, 0)$$

$$p(\mathbf{w}|\boldsymbol{\beta}) \sim \prod_{m=1}^M \mathcal{N}(w_m; 0, \beta^{-1}) \quad \text{and} \quad \beta \sim \Gamma(0, 0)$$

Note that, as “precision” $\alpha_n \rightarrow \infty$, the coefficient x_n is zeroed.

- The *conditional* posterior is (due to Gaussianity) simply

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

- In practice, $\{\boldsymbol{\alpha}, \beta\}$ are estimated using the EM algorithm and then plugged into $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ to approximate the posterior $p(\mathbf{x}|\mathbf{y})$.
- The RVM (also known as “SBL” and “BCS”) is relatively slow.

Other Bayesian Methods

- Bayesian matching pursuits:
 - Greedy methods that use probabilistic support selection.
- Approximate message passing (AMP):
 - Inspired by methods from statistical physics and information theory.
 - Near-optimal in terms of speed and accuracy if \mathbf{A} is large & random.

Alternating Minimization

- also called "cyclic minimization" or "block coordinate descent"
- Cost function to be minimized: $f(\theta)$
- Partition $\theta = \theta_{1:m} = \{\theta_1, \dots, \theta_m\} = \{\theta_{1:k-1}, \theta_k, \theta_{k+1:m}\}$.
- Sweep through the partition. In sweep i :

$$\theta_k^{(i)} = \arg \min_{\theta_k} f(\theta_{1:k-1}^{(i)}, \theta_k, \theta_{k+1:m}^{(i-1)}) \quad k = 1, \dots, m$$

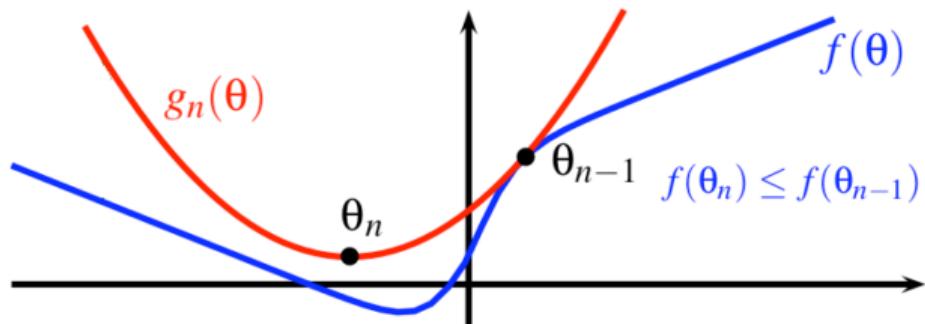
Partitioning done in such a way that these reduced optimization problems are "easy".

- Guaranteed to converge to a local minimum:

$$f(\theta_{1:k}^{(i)}, \theta_{k+1:m}^{(i-1)}) \leq f(\theta_{1:k-1}^{(i)}, \theta_{k:m}^{(i-1)})$$

- Not necessary to do regular sweeps through the partition, can minimize in any order.

Majorization Minimization



- Cost function to be minimized: $f(\theta)$. At iteration $n-1$: have $\theta^{(n-1)}$.
- A **majorizer** at iteration $n-1$ is a function $g_n(\theta)$ such that

$$g_n(\theta^{(n-1)}) = f(\theta^{(n-1)})$$

$$g_n(\theta) \geq f(\theta), \forall \theta$$

- where $g_n(\theta)$ is possibly convex or in any case "easy" to minimize or **just decrease**.
- Then

$$\theta^{(n)} = \arg \min_{\theta} g_n(\theta)$$

$$f(\theta^{(n)}) \leq g_n(\theta^{(n)}) \leq g_n(\theta^{(n-1)}) = f(\theta^{(n-1)})$$

Hence guaranteed to converge to a local minimum.



Expectation Maximization (EM) Algorithm

- Maximum Likelihood estimation: $f(\theta) = -\ln p(\mathbf{y}|\theta)$
- Often resulting from eliminating random \mathbf{x} in $p(\mathbf{y}, \mathbf{x}|\theta)$. θ : deterministic (actual) parameters, \mathbf{x} : random, "unobserved" data, $\{\mathbf{x}, \mathbf{y}\}$: "complete" data.
Estimating θ from $\{\mathbf{x}, \mathbf{y}\}$ easier than from \mathbf{y} alone.
- Majorizer: $g_i(\theta) = f(\theta^{(i-1)}) - E_{\mathbf{x}|\mathbf{y}, \theta^{(i-1)}} \left\{ \ln \left(\frac{p(\mathbf{y}, \mathbf{x}|\theta)}{p(\mathbf{y}, \mathbf{x}|\theta^{(i-1)})} \right) \right\}$, $g_i(\theta^{(i-1)}) = f(\theta^{(i-1)})$

$$\begin{aligned}
 g_i(\theta) &\geq f(\theta^{(i-1)}) - \ln \left(E_{\mathbf{x}|\mathbf{y}, \theta^{(i-1)}} \left\{ \frac{p(\mathbf{y}, \mathbf{x}|\theta)}{p(\mathbf{y}, \mathbf{x}|\theta^{(i-1)})} \right\} \right) \quad \text{Jensen's inequality} \\
 &= f(\theta^{(i-1)}) - \ln \left(E_{\mathbf{x}|\mathbf{y}, \theta^{(i-1)}} \left\{ \frac{p(\mathbf{y}, \mathbf{x}|\theta)}{p(\mathbf{x}|\mathbf{y}, \theta^{(i-1)}) p(\mathbf{y}|\theta^{(i-1)})} \right\} \right) \\
 &= f(\theta^{(i-1)}) - \ln \left(\int p(\mathbf{x}|\mathbf{y}, \theta^{(i-1)}) \frac{p(\mathbf{y}, \mathbf{x}|\theta)}{p(\mathbf{x}|\mathbf{y}, \theta^{(i-1)}) p(\mathbf{y}|\theta^{(i-1)})} d\mathbf{x} \right) \\
 &= f(\theta^{(i-1)}) - \ln \left(\frac{1}{p(\mathbf{y}|\theta^{(i-1)})} \underbrace{\int p(\mathbf{y}, \mathbf{x}|\theta) d\mathbf{x}}_{p(\mathbf{y}|\theta)} \right) = f(\theta^{(i-1)}) - \ln \left(\frac{p(\mathbf{y}|\theta)}{p(\mathbf{y}|\theta^{(i-1)})} \right) = -\ln(p(\mathbf{y}|\theta)) = f(\theta)
 \end{aligned}$$

- EM algorithm: at iteration i :

Expectation step: $\bar{g}_i(\theta) = E_{\mathbf{x}|\mathbf{y}, \theta^{(i-1)}} \{ \ln p(\mathbf{y}, \mathbf{x}|\theta) \}$
 Maximization step: $\theta^{(i)} = \arg \max_{\theta} \bar{g}_i(\theta)$

where $g_i(\theta) = \text{constant} - \bar{g}_i(\theta)$. EM converges to (local) Maximum Likelihood estimate.



Time Varying Sparse State Tracking

Sparse signal \mathbf{x}_t is modeled using an AR(1) process with diagonal correlation coefficient matrix \mathbf{F} .

$$\begin{array}{c} \mathbf{y}_t \\ \vdots \\ N \times 1 \end{array} = \mathbf{A}_t \quad \begin{array}{c} \mathbf{A}_t \\ \vdots \\ N \times M, N \ll M \end{array} + \begin{array}{c} \mathbf{x}_t \\ \vdots \\ N \times 1 \end{array} + \begin{array}{c} \mathbf{v}_t \\ \vdots \\ N \times 1 \end{array}$$

$$\begin{array}{c} \mathbf{x}_t \\ \vdots \\ M \times 1 \end{array} = \mathbf{F} \quad \begin{array}{c} \mathbf{F} \\ \vdots \\ M \times M \end{array} + \begin{array}{c} \mathbf{x}_{t-1} \\ \vdots \\ M \times 1 \end{array} + \begin{array}{c} \mathbf{w}_t \\ \vdots \\ M \times 1 \end{array}$$

Define: $\Xi = \text{diag}(\xi)$, $\mathbf{F} = \text{diag}(\mathbf{f})$.

f_i : correlation coefficient and $x_{i,t} \sim \mathcal{CN}(x_{i,t}; 0, \frac{1}{\xi_i})$. Further, $\mathbf{w}_t \sim \mathcal{CN}(\mathbf{w}_t; \mathbf{0}, \boldsymbol{\Gamma}^{-1} = \Xi^{-1}(\mathbf{I} - \mathbf{FF}^H))$

and $\mathbf{v}_t \sim \mathcal{CN}(\mathbf{v}_t; \mathbf{0}, \gamma^{-1}\mathbf{I})$. **VB leads to Gaussian SAVE-Kalman Filtering (GS-KF).**

Applications: Localization, Adaptive Filtering.

Compressed Sensing Problem

- **Noiseless case:** Given underdetermined \mathbf{y} , \mathbf{A} , the optimization problem is

$$\min_{\mathbf{x}} \|\mathbf{x}\|_0 \text{ subject to } \mathbf{y} = \mathbf{Ax}.$$

Can recover \mathbf{x} and its support for small $N - \|\mathbf{x}\|_0$
(small overdetermination if support were known)

- **Noisy case:**

$$\min_{\mathbf{x}} \|\mathbf{x}\|_0 \text{ subject to } \|\mathbf{y} - \mathbf{Ax}\|_2 \leq \epsilon.$$

- ℓ_0 norm minimization: an NP-hard problem.
- Constrained problem \Rightarrow **Lagrangian, Convex Relaxation** (using p norm, $p > 1$):

$$\min_{\mathbf{x}} \|\mathbf{y} - \mathbf{Ax}\|_2^2 + \lambda \|\mathbf{x}\|_p.$$

Restricted Isometry Property (RIP): $\mathbf{A}^T \mathbf{A}$ sufficiently diagonally dominant

- Most identifiability work considered **noiseless data & exact sparsity**

Sparse Signal Recovery Algorithms

Convex Relaxation based Methods:

- Basis pursuit (ℓ_1 norm)¹.
- LASSO(ℓ_1 norm)²
- Dantzig selector³
- FOCUSS (ℓ_p norm, with $p < 1$).

Greedy Algorithms:

- Matching Pursuit⁴
- Orthogonal Matching Pursuit (OMP)⁵
- CoSaMP⁶

Iterative Methods:

- Iterative Shrinkage and Thresholding Algorithm (ISTA)⁷ or Fast ISTA.
- Approximate Message Passing variants (xAMP)- more details to follow.
- Very recent: Deep learning based methods such as Learned ISTA (LISTA)⁸, Learned AMP (LAMP) and Learned Vector AMP (LVAMP)⁹.

¹Chen, Donoho, Saunders'99, ²Tibshirani'96, ³Candes, Tao'07

⁴Mallat, Zhang'93, ⁵Tropp, Gilbert'07, ⁶Needell, Tropp'09

⁷Daubechies, Defrise, Mol'04, ⁸Gregor, Cun'10, ⁹Borgerding, Schniter, Rangan'17

James-Stein Estimator

- Bayesian interpretation of (possibly overdetermined) Compressed Sensing:

$$\min_x \|y - Ax\|_2^2 - 2\sigma_v^2 \ln p(x)$$

- Stein and James¹⁰ showed that for i.i.d. linear Gaussian model $p(x) = \mathcal{N}(x; \mathbf{0}, \xi^{-1} \mathbf{I})$, it is possible to construct a nonlinear estimate of x with lower MSE than that of ML for all values of the true unknown x .
- A popular design strategy: is to minimize Stein's unbiased risk estimate (SURE), which is an unbiased estimate of the MSE.
- SURE directly approximates the MSE of an estimate from the data, without requiring knowledge of the hyperparameters (ξ), it is an instance of empirical Bayes.
- Stein's landmark discovery lead to the study of biased estimators that outperform minimum variance unbiased estimators (MVU) in terms of MSE, e.g. work by Yonina Eldar¹¹.
- Shrinkage estimators and penalized maximum likelihood (PML) estimators.

¹⁰James, Stein'61

¹¹Eldar'08

Kernel Methods in Automatic Control

- Kernel methods in linear system identification¹² ($\mathbf{y} = \mathbf{Ax} + \mathbf{v}$, $\mathbf{v} \sim \mathcal{N}(\mathbf{v}; \mathbf{0}, \gamma^{-1}\mathbf{I})$).
- Traditional methods: maximum likelihood (ML) or prediction error methods (PEM)
- ML/PEM optimal in the large data limit.
- Questions: Model structure design for ML/PEM. Achieving a good bias-variance trade off.
- Solution: Parameterized Kernel design and hyperparameter estimation. Methods for hyperparameter estimation include cross-validation (CV), empirical Bayes (EB), C_p statistics and Stein's unbiased risk estimate (SURE).
- Regularized least square estimator (\mathbf{P} is symmetric and +ve semidefinite kernel matrix):

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathcal{R}^M} \|\mathbf{y} - \mathbf{Ax}\|^2 + \frac{1}{\gamma} \mathbf{x}^T \mathbf{P}^{-1} \mathbf{x}.$$

- Parameterized family of matrices, $\mathbf{P}(\boldsymbol{\eta})$, where $\boldsymbol{\eta} \in \mathcal{R}^P$. $\boldsymbol{\eta}$ are the hyperparameters.
- Can be interpreted as a constrained form of SBL, with a zero-mean Gaussian prior for \mathbf{x} of which the covariance matrix is a linear combination of some fixed matrices (SBL being a special case with fixed matrices $\mathbf{e}_i \mathbf{e}_i^T$).
- A good overview of Kernel methods, connection with machine learning¹³.

¹²Pillonetto, Nicolao'10

¹³Pillonetto, Dinuzzo, Chen, Nicolao, Ljung'14

Kernel Hyperparameter Estimation

- Empirical Bayes (EB=Type II ML):

$$\widehat{\boldsymbol{\eta}}_{EB} = \arg \min_{\boldsymbol{\eta}} f_{EB}(\boldsymbol{P}(\boldsymbol{\eta})),$$

$$f_{EB}(\boldsymbol{P}(\boldsymbol{\eta})) = \mathbf{y}^T \boldsymbol{Q}^{-1} \mathbf{y} + \ln \det(\boldsymbol{Q}) \text{ with } \boldsymbol{Q} = \mathbf{A} \boldsymbol{P} \mathbf{A}^T + \frac{1}{\gamma} \mathbf{I}_N.$$

- Stein's Unbiased Risk Estimator (SURE) method:

- SURE: MSE of signal reconstruction ($MSE_x(\boldsymbol{P}) = E(\|\hat{\mathbf{x}} - \mathbf{x}\|^2)$):

$$SURE_x : \widehat{\boldsymbol{\eta}}_{Sx} = \arg \min_{\boldsymbol{\eta}} f_{Sx}(\boldsymbol{P}(\boldsymbol{\eta})), \text{ with}$$

$$f_{Sx}(\boldsymbol{P}(\boldsymbol{\eta})) = \frac{1}{\gamma^2} \mathbf{y}^T \boldsymbol{Q}^{-T} \mathbf{A} (\mathbf{A}^T \mathbf{A})^{-2} \mathbf{A}^T \boldsymbol{Q}^{-1} \mathbf{y} + \frac{1}{\gamma} \text{tr}\{2\boldsymbol{R}^{-1} - (\mathbf{A}^T \mathbf{A})^{-1}\},$$

$$\boldsymbol{R} = \mathbf{A}^T \mathbf{A} + \frac{1}{\gamma} \boldsymbol{P}^{-1}.$$

- The SURE estimator converge to the best possible hyperparameter in terms of MSE in the asymptotic limit, “asymptotically consistent”.
- EB estimator converges to another best hyperparameter minimizing the expectation of the EB estimation criterion (EEB).
- Convergence of EB is faster than that of the SURE estimator.

¹³Mu, Chen, Ljung'18

Outline

1 Parametric Signal Models

2 Compressed Sensing

3 SBL

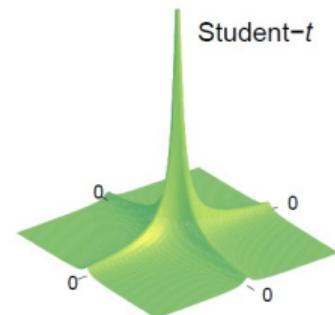
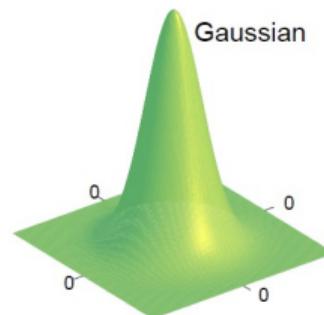
Sparse Bayesian Learning - SBL

- Bayesian Compressed Sensing: 2-layer hierarchical prior for \mathbf{x} as in ¹⁴, inducing sparsity for \mathbf{x} (conjugate priors: posterior pdf of same family as prior pdf) :

$$p_{\mathbf{x}}(x_{i,t}|\xi_i) = \mathcal{N}(x_{i,t}; 0, \xi_i^{-1}), \quad p(\xi_i|a, b) = \Gamma^{-1}(a)b^a \xi_i^{a-1} e^{-b\xi_i}$$

⇒ sparsifying Student-t marginal

$$p_{\mathbf{x}}(x_{i,t}) = \frac{b^a \Gamma(a + \frac{1}{2})}{(2\pi)^{\frac{1}{2}} \Gamma(a)} (b + x_{i,t}^2/2)^{-(a + \frac{1}{2})}$$



- Sparsification of the Innovation Sequence: we apply the (Gamma) prior not to the precision of the state \mathbf{x} but of its innovation \mathbf{w} , allowing to sparsify at the same time the components of \mathbf{x} AND their variation in time (innovation).
- The inverse of the noise variance γ is also assumed to have a Gamma prior,
 $p_{\gamma}(\gamma|c, d) = \Gamma^{-1}(c)d^c \gamma^{c-1} e^{-d\gamma}$.

¹⁴Tipping'01

Original SBL Algorithm (Type II ML)

- Original SBL¹⁵, for a fixed estimate of the **hyperparameters** $(\hat{\xi}, \hat{\gamma})$, the posterior of x is Gaussian, i.e.

$$p_x(x|y, \hat{\xi}, \hat{\gamma}) = \mathcal{N}(x; \hat{x}, \Sigma_L)$$

leading to the **(Linear) MMSE estimate for x**

$$\begin{aligned} \hat{x} &= \hat{\gamma}(\hat{\gamma}\mathbf{A}^T\mathbf{A} + \hat{\Xi})^{-1}\mathbf{A}^Ty, \\ \Sigma_L &= (\hat{\gamma}\mathbf{A}^T\mathbf{A} + \hat{\Xi})^{-1}. \end{aligned} \quad (1)$$

- The **hyperparameters** are estimated from the likelihood function by marginalizing over the sparse coefficients x , the marginalized likelihood being denoted as $p_y(y|\xi, \gamma)$. ξ, γ are estimated by maximizing $p_y(y|\xi, \gamma)$ and this procedure is called as Type-II ML. **Type-II ML is solved using EM**, which leads to the following updates for the hyperparameters.

$$\hat{\xi}_i = \frac{a + \frac{1}{2}}{\left(\frac{\langle x_i^2 \rangle}{2} + b \right)}, \text{ where } \langle x_i^2 \rangle = \hat{x}_i^2 + \sigma_i^2. \quad \langle \gamma \rangle = \frac{c + \frac{N}{2}}{\left(\frac{\langle \|y - Ax\|^2 \rangle}{2} + d \right)},$$

$$\text{where, } \langle \|y - Ax\|^2 \rangle = \|y\|^2 - 2y^T\mathbf{A}\hat{x} + \text{tr}(\mathbf{A}^T\mathbf{A}(\hat{x}\hat{x}^T + \Sigma)), \\ \Sigma = \text{diag}(\Sigma_L) = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_M^2), \quad \hat{x} = [\hat{x}_1, \hat{x}_2, \dots, \hat{x}_M]^T.$$

¹⁵Tipping'01, Wipf,Rao'04

Type I vs Type II ML

- Type I \Rightarrow standard MAP estimation (involves integrating out the hyperparameters)

$$\hat{\mathbf{x}} = \arg \max_{\mathbf{x}} [\log p_{\mathbf{y}}(\mathbf{y}|\mathbf{x}) + p_{\mathbf{x}}(\mathbf{x})].$$

- Type II \Rightarrow hyperparameters ($\Psi = \{\xi, \gamma\}$) are estimated using an evidence maximization approach

$$\hat{\Psi} = \arg \max_{\Psi} p_{\Psi}(\Psi|\mathbf{y}) = \arg \max_{\Psi} p_{\Psi}(\Psi) \int p_{\mathbf{y}}(\mathbf{y}|\Psi) = \arg \max_{\Psi} p_{\Psi}(\Psi) \int p_{\mathbf{y}}(\mathbf{y}|\mathbf{x}, \gamma) p_{\mathbf{x}}(\mathbf{x}|\xi) d\mathbf{x}.$$

- Why Type II is better than Type I? ¹⁶ In the evidence maximization framework instead of looking for the mode of the true posterior $p_{\mathbf{x}}(\mathbf{x}|\mathbf{y})$, the true posterior is approximated as $p_{\mathbf{x}}(\mathbf{x}|\mathbf{y}; \hat{\Psi})$, where $\hat{\Psi}$ is obtained by maximizing the true posterior mass over the subspaces spanned by the non zero indexes.
- Type I methods seek the mode of the true posterior and use that as the point estimate of the desired coefficients. Hence, if the true posterior distribution has a skewed peak, then the Type I estimate (Mode) is not a good representative of the whole posterior.

¹⁶Giri, Rao'16

Variational Bayesian (VB) Inference

- The computation of the posterior distribution of the parameters is usually intractable. As in SAGE, **SAVE is simply VB with partitioning of the unknowns at the scalar level**. Define $\theta = \{x, \xi, \gamma\}$, θ_i represents each scalar and $\theta_{\bar{i}}$ denotes θ excluding θ_i .

$$q(\theta) = q_\gamma(\gamma) \prod_{i=1}^M q_{x_i}(x_i) \prod_{i=1}^M q_{\xi_i}(\xi_i).$$

- VB compute the factors q by **minimizing the Kullback-Leibler distance** between the true posterior distribution $p(\theta|y)$ and the $q(\theta)$. From ¹⁷,

$$KLD_{VB} = D_{KL}(q(\theta)||p(\theta|y)) = \int q(\theta) \ln \frac{q(\theta)}{p(\theta|y)} d\theta.$$

- The KL divergence minimization is equivalent to **maximizing the evidence lower bound (ELBO)**¹⁸.

$$\ln p(y) = L(q) + KLD_{VB} = -D_{KL}(q(\theta)||p(\theta, y)) + D_{KL}(q(\theta)||p(\theta|y)), \text{ where,}$$

$\ln p(y)$ is the evidence, and $\min KLD_{VB}$ becomes equivalent to $\max L(q)$, the ELBO.

- We get for the element-wise VB recursions: **(Expectation Maximization (EM))** = special case:

$$\ln(q_i(\theta_i)) = \langle \ln p(y, \theta) \rangle_{\theta_{\bar{i}}} + c_i,$$

θ_s random, hidden
 θ_d deterministic)

¹⁸Beal'03, ¹⁹Tzikas, Likas, Galatsanos'08

Low Complexity-Space Alternating Variational Estimation (SAVE)

- Mean Field (MF) approximation: VB partitioned to scalar level (MF vs VB // SAGE vs EM), results in a SBL algorithm **without any matrix inversions**.
- The resulting **alternating optimization of the posteriors for each scalar in θ** leads to

$$\ln(q_i(\theta_i)) = \langle \ln p(\mathbf{y}, \theta) \rangle_{k \neq i} + c_i,$$

$$p(\mathbf{y}, \theta) = p_{\mathbf{y}}(\mathbf{y}|\mathbf{x}, \xi, \gamma)p_{\mathbf{x}}(\mathbf{x}|\xi)p_{\xi}(\xi)p_{\gamma}(\gamma).$$

where $\theta = \{\mathbf{x}, \xi, \gamma\}$ and θ_i represents each scalar in θ .

$$\begin{aligned} \ln p(\mathbf{y}, \theta) &= \frac{N}{2} \ln \gamma - \frac{\gamma}{2} \|\mathbf{y} - \mathbf{Ax}\|^2 + \sum_{i=1}^M \left(\frac{1}{2} \ln \xi_i - \frac{\xi_i}{2} x_i^2 \right) + \\ &\quad \sum_{i=1}^M ((a-1) \ln \xi_i + a \ln b - b \xi_i) + (c-1) \ln \gamma + c \ln d - d \gamma + \text{constants}. \end{aligned}$$

- Gaussian approximate posterior for x_i :**

$$\begin{aligned} \ln q_{x_i}(x_i) &= -\frac{\gamma}{2} \left\{ \langle \|\mathbf{y} - \mathbf{A}_{\bar{i}} \mathbf{x}_{\bar{i}}\|^2 \rangle - (\mathbf{y} - \mathbf{A}_{\bar{i}} \langle \mathbf{x}_{\bar{i}} \rangle)^T \mathbf{A}_i \mathbf{x}_i - \right. \\ &\quad \left. x_i \mathbf{A}_i^T (\mathbf{y} - \mathbf{A}_{\bar{i}} \langle \mathbf{x}_{\bar{i}} \rangle) + \|\mathbf{A}_i\|^2 x_i^2 \right\} - \frac{\xi_i}{2} x_i^2 + c_{x_i} = -\frac{1}{2\sigma_i^2} (x_i - \hat{x}_i)^2 + c'_{x_i}. \end{aligned}$$

SAVE Iterations Continued...

- The SAVE iterations for \mathbf{x} get obtained as

$$\sigma_i^2 = \frac{1}{\gamma \|\mathbf{A}_i\|^2 + \xi_i}, \quad \hat{\mathbf{x}}_i = \sigma_i^2 \mathbf{A}_i^T (\mathbf{y} - \mathbf{A}_i \hat{\mathbf{x}}_i) \gamma.$$

- Hyperparameter estimates: same as EM iterations. Gamma posterior for ξ_i and γ .
- No matrix inversions.
- Update of all the variables, $\mathbf{x}, \xi_i, \gamma$, requires simple addition and multiplication operations
- $\mathbf{y}^T \mathbf{A}$, $\mathbf{A}^T \mathbf{A}$ and $\|\mathbf{y}\|^2$ can be precomputed, so only need to be computed once.

- Remarks:** From LMMSE expression in (1), i^{th} row of $\gamma \mathbf{A}^T \mathbf{A} \hat{\mathbf{x}} + \Xi \hat{\mathbf{x}} = \gamma \mathbf{A}^T \mathbf{y}$:

$$\gamma \mathbf{A}_i^T \mathbf{A} \hat{\mathbf{x}} + \xi_i \hat{\mathbf{x}}_i = \gamma \mathbf{A}_i^T \mathbf{y} \text{ or } (\gamma \|\mathbf{A}_i\|^2 + \xi_i) \hat{\mathbf{x}}_i = \gamma \mathbf{A}_i^T (\mathbf{y} - \mathbf{A}_i \hat{\mathbf{x}}_i)$$

Hence **SAVE does linear PIC iterations to compute the LMMSE estimate.**

- However, for the posterior variances : $\sigma_i^2 = ((\Sigma_L^{-1})_{i,i})^{-1} \leq (\Sigma_L)_{ii}$, i with equality only for diagonal Σ_L

Convergence of SAVE

Theorem 1

The convergence condition for the sparse coefficients x_i of the SAVE algorithm²⁰ can be written as $\rho(\mathbf{D}^{-1}\mathbf{H}) < 1$, where $\mathbf{D} = \text{diag}(\hat{\gamma}\mathbf{A}^T\mathbf{A} + \hat{\Xi})$, $\mathbf{H} = \text{offdiag}(\hat{\gamma}\mathbf{A}^T\mathbf{A})$. $\rho(\cdot)$ denotes the spectral radius. Moreover, if SAVE converges, assuming the estimate of hyperparameters are consistent, the posterior mean (point estimate) always converges to the exact value (LMMSE). However, the predicted posterior variance is quite suboptimal.

Remark: To fix the convergence of SAVE (when $\rho(\mathbf{D}^{-1}\mathbf{H}) > 1$), we can use the diagonal loading method²¹. The modified iterations (with a diagonal loading factor matrix Λ) can be written as,

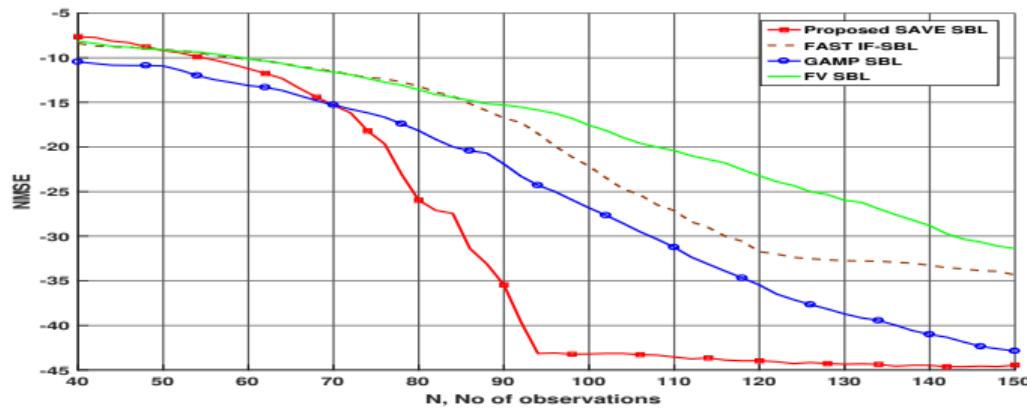
$$\begin{aligned} (\mathbf{D} + \tilde{\Xi})\mathbf{x}^{(t+1)} &= -(\mathbf{H} - \tilde{\Xi})\mathbf{x}^{(t)} + \hat{\gamma}\mathbf{A}^T\mathbf{y}, \implies \\ \mathbf{x}^{(t+1)} &= -(\mathbf{D} + \tilde{\Xi})^{-1}(\mathbf{H} - \tilde{\Xi})\mathbf{x}^{(t)} + (\mathbf{D} + \tilde{\Xi})^{-1}\hat{\gamma}\mathbf{A}^T\mathbf{y}. \end{aligned}$$

The convergence condition gets modified as $\rho((\mathbf{D} + \tilde{\Xi})^{-1}(\mathbf{H} - \tilde{\Xi})) < 1$. Another point worth noting here is that, if the power delay profile Ξ is also estimated using MF, $\hat{\gamma}\text{diag}(\mathbf{A}^T\mathbf{A}) + \hat{\Xi}$, where $\tilde{\Xi} = \Xi + \Xi$, with $\tilde{\Xi} > 0$. In this case, $\tilde{\Xi}$ may represent an automatic correction factor (diagonal loading) to force convergence of SAVE for cases where $\rho(\mathbf{D}^{-1}\mathbf{H}) > 1$.

²⁰Thomas, Slock'18

²¹Johnson, Bickson, Dolev'09

NMSE Results



NMSE vs the number of observations M ($N = 200, K = 40$).

- For sufficient amount of data, **SAVE has significantly lower MSE than the other fast algorithms.**
- Why? The resulting problem of alternating optimization of x and the hyperparameters ξ and γ appears to be characterized by many local optima. Apparently, **the component-wise VB approach appears to allow to avoid a lot of bad local optima, explaining the better performance, apart from lower complexity.**
- At very low amount of data, suboptimal approaches such as AMP which don't introduce individual hyper parameters per x component and assume that the x_i behave i.i.d. behave better because of the lower number of hyper parameters to be estimated.

An Overview of Fast SBL Algorithms

- Fast SBL using Type II ML by Tipping²²: greedy approach handling one x_i at a time, plus replacing precisions by their convergence values, leading to pruning of the small x_i components, i.e. explicit sparsity.
- Fast SBL using VB by Shutin et. al.²³: Shutin uses VB while Tipping is Type II ML as in original SBL. They do both replace precisions by their convergence values. Shutin also added some extra view points in terms of the pruning condition being interpreted as relating between sparsity properties of SBL and a measure of SNR. Main message of the both being faster convergence compared to original SBL, not much reduction in per iteration complexity.
- BP-SBL²⁴: In SBL, with fixed hyperparameters, MAP or MMSE estimate (follows from the Gaussian posterior) of x can be efficiently computed using belief propagation (BP), since all the messages involved are Gaussian (without any approx.).
- Inverse Free SBL (IF-SBL)²⁵: Optimization using a relaxed ELBO.
- Hyperparameter free sparse estimation²⁶: Does not require hyperparameter tuning compared to SBL. Uses covariance matching, equivalent to square root LASSO.

²¹Tipping, Faul'03, ²²Shutin, Buchgraber, Kulkarni, Poor'11

²³Tan, Li'10, ²⁴Duan, Yang, Fang, Li'17

²⁵Zachariah, Stoica'15

Complexity Comparisons-SBL Algorithms

| Algorithm | Complexity per Iteration | Convergence (No of iterations) | Sparsity | Optimization function | Local Optimum |
|---|--|--------------------------------|---|---|---|
| Type I | $O(M^3)$ | | Exact sparsity | Type I ML (Depending upon the prior used, type 1 ML corresponds to LASSO/Re-weighted l1/l2 min. problems) | |
| Type II SBL | $O(M^3)$ | | Exact sparsity (α_i converges to ∞) | Type II ML solved using EM | |
| Fast SBL using Type II ML (Tipping,Faul'03) (Focus more on Convergence speed) | $O(L^3)$, $L \leq M$ | $\ll L$ | Exact sparsity (Using an entry dependent thresholding condition which follows from the computation of stationary point of α_i) | Type II ML (stationary points of α_i are computed to accelerate convergence) | Convergence to a local optimum. |
| Fast SBL (using VB) by Shutin (Focus more on Convergence speed) | $O(L^3)$, $L \leq M$ | $\ll L$ | Exact sparsity (Using a pruning condition similar as in Tipping's) | Maximization of ELBO in VB | Convergence to a local optimum of ELBO (Mean field free energy) |
| Hyperparameter free SBL (Zachariah, Stoica'15) | $O(M^2)$ | $\ll M$ | The final objective function is a weighted square root LASSO. So the sum of l2 norm of (y and Ax) and weighted l1 norm of x which promotes sparsity here. | MMSE estimator for x with Covariance matching for PDP, finally giving rise to an objective function which can be interpreted as weighted square root LASSO. | Convergence to a local optimum |
| BP-SBL (Tan, Li'10) | $O(MN)$ (Similar complexity as xAMP, see matrix form of the BP-SBL in the upcoming slides) | $\log(MN)$ | Does not give exact sparsity | Posterior of x computed using BP and EM for hyper-parameters | Convergence to local optimum of Bethe Free Energy (BFE) |
| GAMP-SBL (Shoukairi, Schniter, Rao'18) | $O(MN)$ | $\ll M$ | Does not give exact sparsity | Using GAMP for posterior of x, EM for hyperparameters | Convergence to local optimum of LSL-BFE |
| SAVE | $O(MN)$ | $\ll M$ | Does not give exact sparsity | Maximization of ELBO in VB | Convergence to a local optimum of ELBO |
| Inverse Free SBL (Duan, Yang, Fang, Li'17) | $O(MN)$ | $\ll M$ (similar to GAMP SBL) | Does not give Exact sparsity | Maximization of an approximate ELBO in VB | Convergence to a local opt in the approx ELBO |

References I

-  M. Bayati, A. Montanari, "The Dynamics of Message Passing on Dense Graphs, with Applications to Compressed Sensing," in *IEEE Trans. on Info. Theo.*, Feb. 2011.
-  M. J. Beal, "Variational Algorithms for Approximate Bayesian Inference," in *Thesis, Univeristy of Cambridge, UK*, May 2003.
-  M. Borgerding, P. Schniter, S. Rangan, "AMP-Inspired Deep Networks for Sparse Linear Inverse Problems," *IEEE Trans. on Sig. Process.*, Aug. 2017.
-  B. Çakmak, M. Opper, "Expectation Propagation for Approximate Inference: Free Probability Framework," in *IEEE Inter. Sympo. On Info. Theo. (ISIT)*, 2018.
-  B. Çakmak, O. Winther, B. H. Fleury, "S-AMP: Approximate Message Passing for General Matrix Ensembles," in *IEEE Intl. Sympo. Info. Theo.*, 2014.
-  E. Candes, T. Tao, "The Dantzig selector: Statistical estimation when p is much larger than n ," *The annals of Statistics*, 2007.
-  S. S. Chen, D. L. Donoho, M. A. Saunders, "Atomic decomposition by Basis Pursuit," *SIAM J. Sci. Comput.*, 1999.
-  R. Couillet, J. Hoydis, M. Debbah", "Random Beamforming over Quasi-Static and Fading Channels: A Deterministic Equivalent Approach," in *IEEE Trans. On Info. Theo.*, 2012.
-  I. Daubechies, M. Defrise, C. D. Mol, "An iterative thresholding algorithm for linear inverse problems with a sparsity constraint..," *Comm. on Pure and Applied Mathematics*, 2004.
-  J. Du, S. Ma, Y-C. Wu, S. Kar, J. M. F. Moura, "Convergence Analysis of Distributed Inference with Vector-Valued Gaussian Belief Propagation," in *Jrnl. of Mach. Learn. Res.*, Apr. 2018.
-  H. Duan, L. Yang, J. Fang, H. Li, "Fast Inverse-Free Sparse Bayesian Learning via Relaxed Evidence Lower Bound Maximization," in *IEEE Sig. Process. Lett.*, Jun. 2017.
-  Y. C. Eldar "Rethinking Biased Estimation: Improving Maximum Likelihood and the CramérRao Bound," *Found. and Tren. in Sig. Process.*, 2008.

References II

-  S. Fortunati, F. Gini, M. S. Greco, C. D. Richmond, "Performance Bounds for Parameter Estimation under Misspecified Models: Fundamental Findings and Applications," in *IEEE Sig. Proc. Mag.*, Nov. 2017.
-  A. E. Gelfand, S. K. Sahu, "Identifiability, Improper Priors, and Gibbs Sampling for Generalized Linear Models," in *Journ. of the Americ. Stat. Assoc.*, Mar. 1999.
-  R. Giri, B. Rao, "Type I and Type II Bayesian Methods for Sparse Signal Recovery Using Scale Mixtures," in *IEEE Trans. Sig. Process.*, July 2016.
-  K. Gregor, Y. LeCun, "Learning Fast Approximations of Sparse Coding," *Intl. Conf. on Mach. Learn.*, 2010.
-  W. James, C. Stein "Estimation with quadratic loss," *Proc. 4th Berkeley Symp. Mathematical Statistics Probability*, 1961.
-  J. K. Johnson, D. Bickson, D. Dolev, "Fixing Convergence of Gaussian Belief Propagation," in *IEEE Intl. Symp. on Info. Theo.*, 2009.
-  M. Luo, Q. Guo, D. Huang, J. Xi, "Sparse Bayesian Learning using Approximate Message Passing with Unitary Transformation," in *IEEE VTS Asia Pac. Wire. Commun. Symp., APWCS*, Aug. 2019.
-  D. M. Malioutov, J. K. Johnson, A. S. Willsky, "Walk-Sums and Belief Propagation in Gaussian Graphical Models," in *Jrnl. of Mach. Learn. Res.*, Oct. 2006.
-  S. Mallat, Z. Zhang, "Matching Pursuits with time-frequency dictionaries," *IEEE Trans. Sig. Process.*, 1993.
-  J. Ma, L. Ping "Orthogonal AMP," in *IEEE Access*, Mar. 2017.
-  T. P. Minka, "Expectation propagation for approximate Bayesian inference , " in *Proc. of Conf. on Uncert. in Art. Intell. (UAI)*, 2001.
-  B. Mu, T. Chen, L. Ljung, "On asymptotic properties of hyperparameter estimators for kernel-based regularization methods," in *Automatica*, May. 2018.
-  D. Needell, J. A. Tropp, "CoSaMP: Iterative signal recovery from incomplete and inaccurate samples., " *Applied and computational harmonic analysis ,* 2009.



References III

-  G. Pillonetto, F. Dinuzzo, T. Chen, G. D Nicolao, L. Ljung, "Kernel methods in system identification, machine learning and function estimation: A survey," *Automatica*, Feb. 2014.
-  G. Pillonetto, G. D. Nicolao, "A new kernel-based approach for linear system identification," *Automatica*, 2010.
-  S. Rangan, "Generalized Approximate Message Passing for Estimation with Random Linear Mixing," in *IEEE Intl. Sympo. Info. Theo.*, 2011.
-  S. Rangan, A. K. Fletcher, P. Schniter, U. S. Kamilov, "Inference for Generalized Linear Models via Alternating Directions and Bethe Free Energy Minimization," in *IEEE Trans. Inf. Theory*, Jan. 2017.
-  S. Rangan, P. Schniter, A. K. Fletcher, "Vector Approximate Message Passing," in *IEEE Trans. Inf. Theory*, Oct. 2019.
-  C. D. Richmond and L. L. Horowitz, "Parameter bounds on estimation accuracy under model misspecification," in *IEEE Trans. on Sig. Process.*, May. 2015.
-  E. Riegler, G. E. Kirkelund, C. N. Manchón, B. H. Fleury, "Merging Belief Propagation and the Mean Field Approximation: a Free Energy Approach," in *IEEE Trans. on Info. Theo.*, Jan. 2013.
-  P. Rusmevichtong, B. Van Roy, "An analysis of belief propagation on the turbo decoding graph with Gaussian densities," *IEEE Trans. Inf. Theory*, 2001.
-  M. Al-Shoukairi, P. Schniter, B. D. Rao, "GAMP-Based Low Complexity Sparse Bayesian Learning Algorithm," in *IEEE Trans. on Sig. Process.*, Jan. 2018.
-  D. Shutin, T. Buchgraber, S. R. Kulkarni, H. V. Poor, "Fast Variational Sparse Bayesian Learning With Automatic Relevance Determination for Superimposed Signals," in *IEEE Trans. Sig. Process.*, Dec. 2011.
-  V. Šmídl, A. Quinn "The Variational Bayes Method in Signal Processing," in *Springer Series on Sig. and Comm. Tech.*, 2005.
-  K. Takeuchi, "Rigorous Dynamics of Expectation- Propagation-Based Signal Recovery from Unitarily Invariant Measurements," in *IEEE Trans. Inf. Theory*, Jan. 2020.
-  X. Tan, J. Li, "Computationally Efficient Sparse Bayesian Learning via Belief Propagation," in *IEEE Trans. on Sig. Proc.*, Apr. 2010.

References IV

-  C. K. Thomas, K. Gopala, D. Slock, "Sparse Bayesian learning for a bilinear calibration model and mismatched CRB," in *EUSIPCO*, 2019.
-  C. K. Thomas, D. Slock, "SAVE - space alternating variational estimation for sparse Bayesian learning," in *IEEE Data Sci. Wkshp.*, Jun. 2018.
-  C. K. Thomas, D. Slock, "Space Alternating Variational Estimation and Kronecker Structured Dictionary Learning," in *IEEE ICASSP*, 2019.
-  C. K. Thomas and D. Slock, "Convergence Analysis Of Sparse Bayesian Learning under Approximate Inference Techniques," in *Asilomar Conf. on Sig., Sys., and Comp.*, Nov. 2019.
-  C. K. Thomas and D. Slock, "Low Complexity Static and Dynamic Sparse Bayesian Learning Combining BP, VB and EP Message Passing," in *IEEE Asilomar*, Nov. 2019.
-  R. Tibshirani, "Regression shrinkage and selection via the LASSO," *J. Roy. Statist. Soc. Series B (Methodol.)*, 1996.
-  M. E. Tipping, "Sparse Bayesian Learning and the Relevance Vector Machine," *J. Mach. Learn. Res.*, 2001.
-  M. E. Tipping, A. Faul, "Fast marginal likelihood maximisation for sparse Bayesian models," in *AISTATS*, 2003.
-  J. A. Tropp, A. C. Gilbert, "Signal Recovery From Random Measurements Via Orthogonal Matching Pursuit," *IEEE Trans. Info. Theo.*, 2007.
-  D. G. Tzikas, A. C. Likas, N. P. Galatsanos, "The variational approximation for Bayesian inference," in *IEEE Sig. Process. Mag.*, Nov. 2008.
-  S. Wagner, R. Couillet, M. Debbah, D. Slock, "Large System Analysis of Linear Precoding in MISO Broadcast Channels with Limited Feedback," in *IEEE Trans. Inf. Theory*, July. 2012.
-  Y. Weiss, W. T. Freeman, "Correctness of belief propagation in Gaussian graphical models of arbitrary topology," *NIPS*, 2000.
-  D. P. Wipf, B. D. Rao, "Sparse Bayesian Learning for Basis Selection," *IEEE Trans. Sig. Process.*, Aug 2004.
-  J. S. Yedidia, W. T. Freeman, Y. Weiss, "Constructing free-energy approximations and generalized belief propagation algorithms," in *IEEE Trans. on Info. Theo.*, Jun. 2005.
-  D. Zachariah, P. Stoica, "Online Hyperparameter-Free Sparse Estimation Method," in *IEEE Trans. on Sig. Proc.*, July. 2015.